

Visual feature representation from sparse to deep

Qingshan Liu

Nanjing University of Information Science & Technology

28. 4. 2017

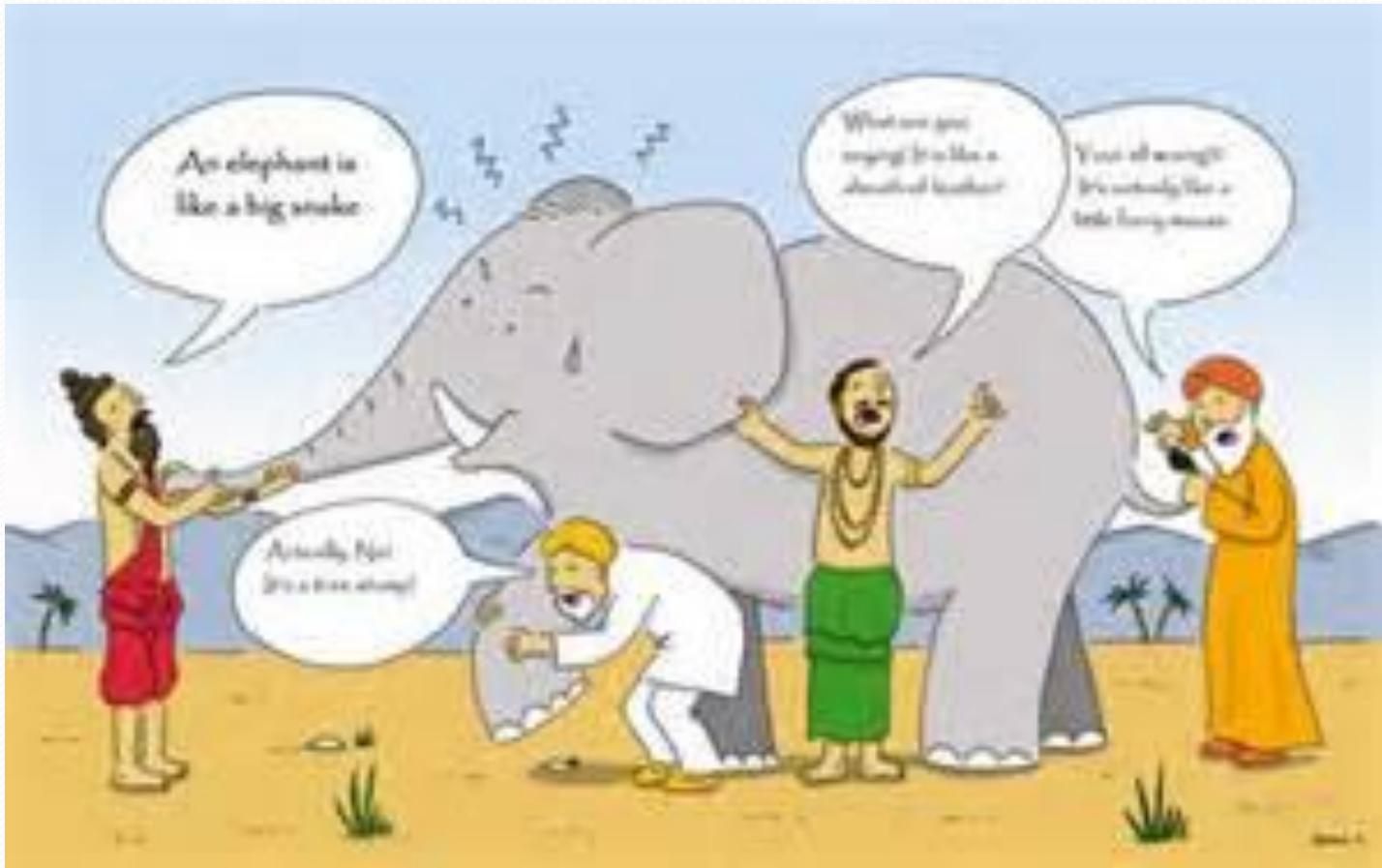
郭晶晶
旺夫脸
的六大特征



南京信息工程大学



江苏省大数据分析技术重点实验室
Jiangsu Key Laboratory of Big Data Analysis Technology



南京信息工程大学



江苏省大数据分析技术重点实验室
Jiangsu Key Laboratory of Big Data Analysis Technology

这是一个及其美好的时代，这是一个极具挑战的时代



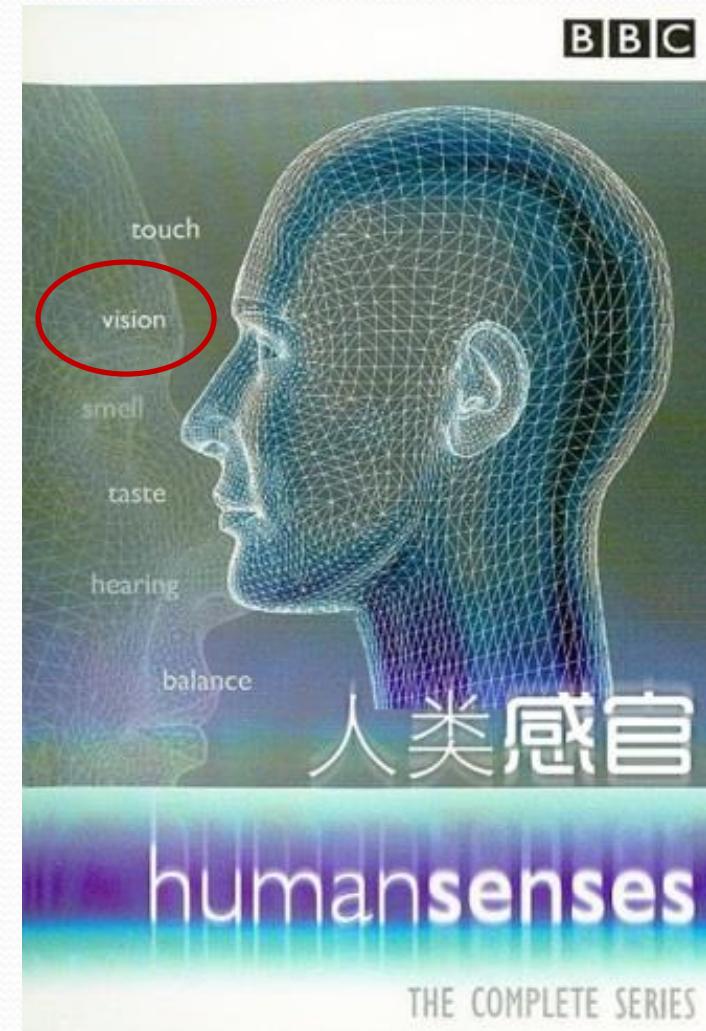
2017年3月5日 十二届全国人大五次会议

李克强总理在《政府工作报告》中指出，要加快培育壮大新兴产业，全面实施包括人工智能在内的战略性新兴产业发展规划。**“人工智能”首次被写入政府工作报告。**

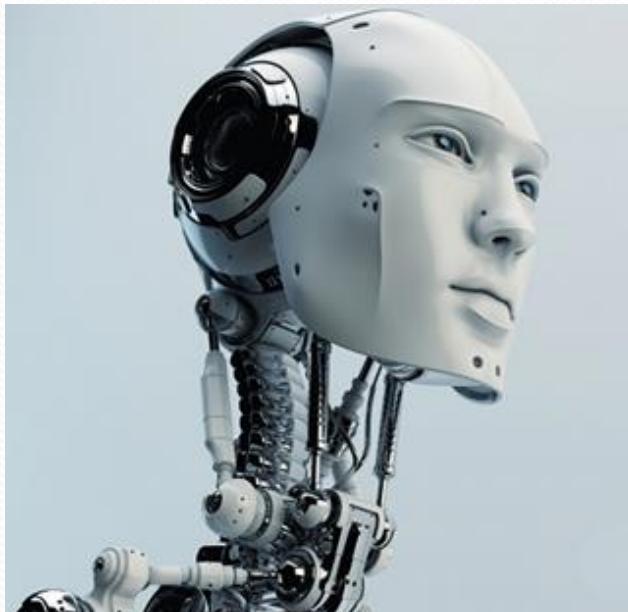
人工智能离不开机器视觉



眼睛是心灵的窗户
眼睛是人类感知外部环境最重要的器官



■ Purpose : Make machine have eyes



Low level features

Semantic gap



Image semantics



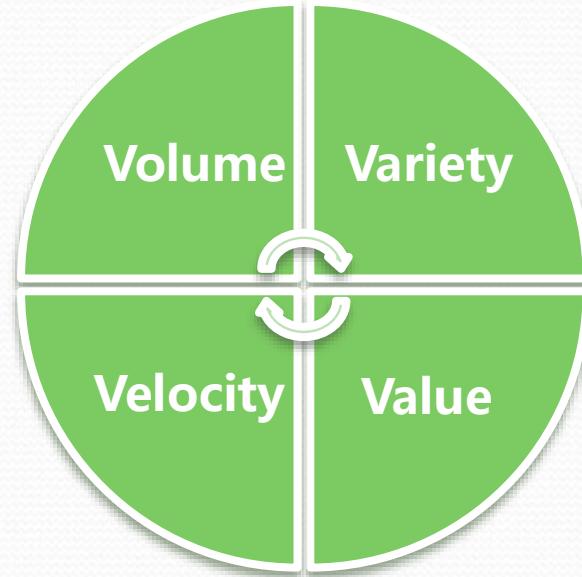
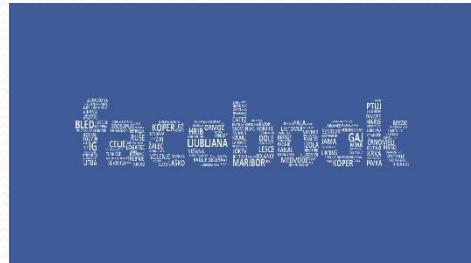
江苏省大数据分析技术重点实验室
Jiangsu Key Laboratory of Big Data Analysis Technology



南京信息工程大学

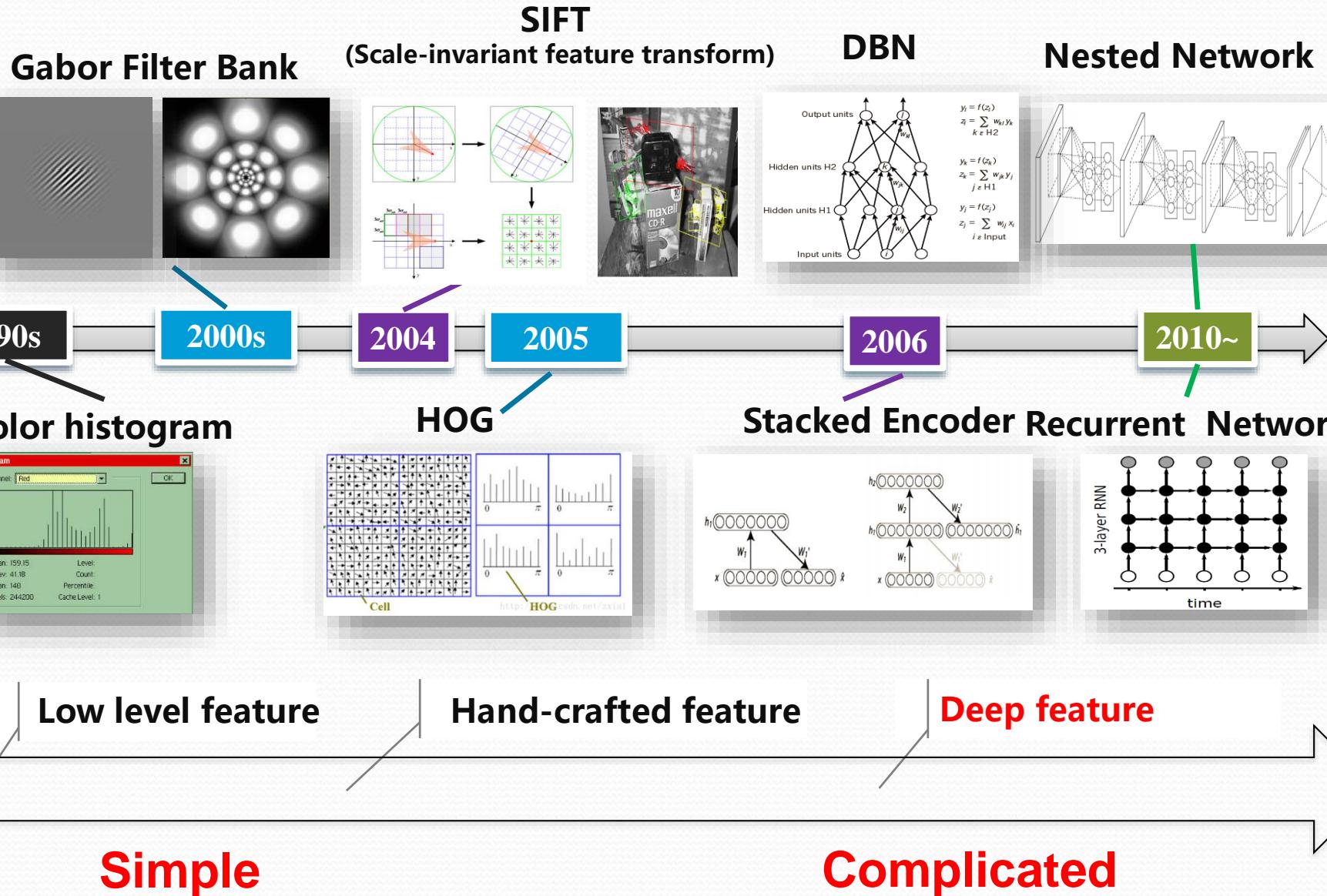
Challenges

- There have 100 million surveillance cameras distributed in the word, which will produce **2.3 ZB (10^{21})** video data
- Youtube will increase over **72 hours** video data in each minute
- Face book has over **300 billion** images
-



Data driven high complexity issue

Visual Feature Representation



Outline

- **Sparse based feature representation**
- **Hypergraph-based feature representation**
- **Deep feature representation**

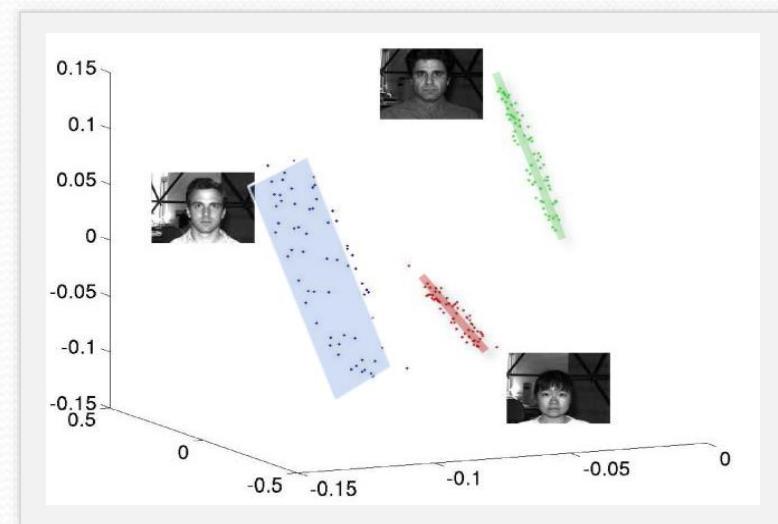
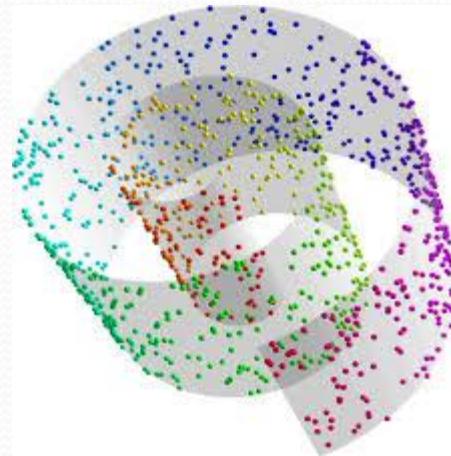


南京信息工
程大学



江苏省大数据分析技术重点实验室
Jiangsu Key Laboratory of Big Data Analysis Technology

David L. Donoho, High-dimensional data analysis: The **curses** and **blessings** of dimensionality. *Aide-Memoire of a Lecture at (2000)*



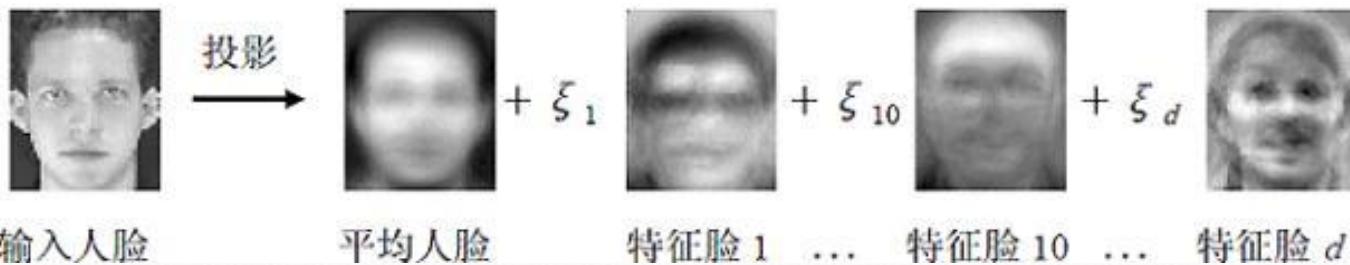
How to learn the low dimensional feature representation

Subspace Learning

- Learn a low dimensional subspace projection to handle the high-dimensional data

$$y = A^T x$$

$$x \in R^D, \quad A \in R^{D \times d}, \quad y \in d, \quad d < D.$$



Subspace Learning

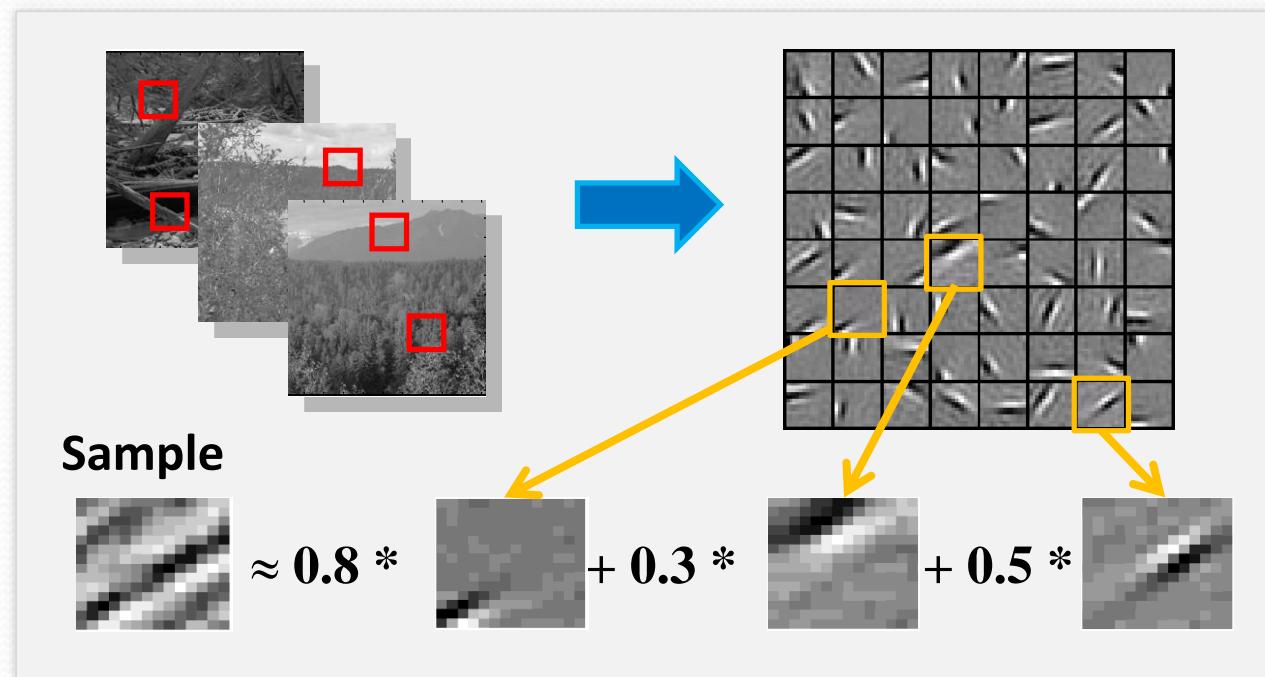
- **Linear subspace:** A is a linear transformation
for example: PCA, LDA,...
- **Kernel based nonlinear subspace:** combining
the nonlinear kernel trick with linear subspace
for example: KPCA, KLDA,...
- **Manifold subspace**
for example: LLE, ISOMap,...

Sparse feature representation

Simple



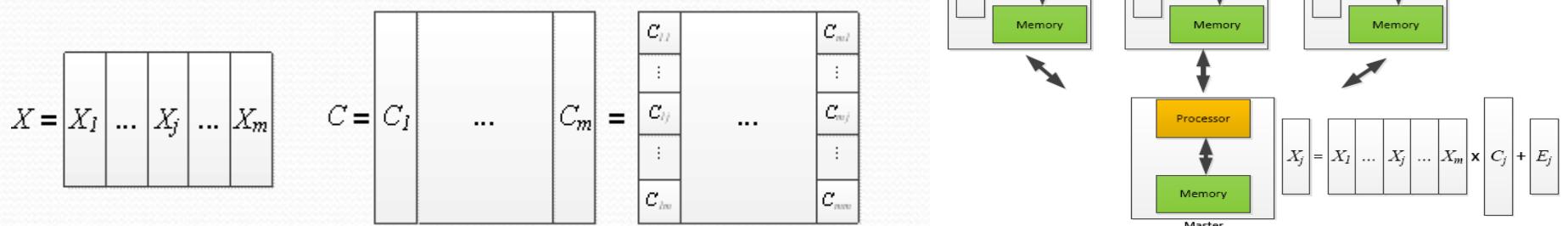
Reliable



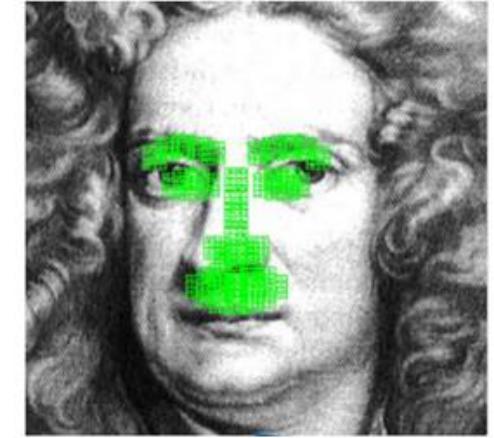
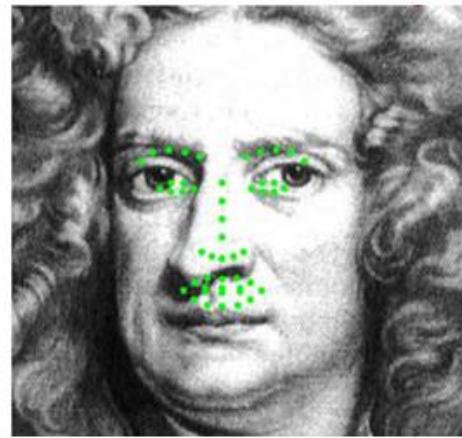
$$\min_{A, X} \|X - AZ\|_2^2 + \lambda \|Z\|_1$$

Sparse Representation

- Learning Discriminative Dictionary for Group Sparse Representation (**IEEE T-IP 2014**)
- Newton Greedy Pursuit: a Quadratic Approximation Method for Sparsity-Constrained Optimization, (**CVPR 2014**).
- Decentralized Robust Subspace Clustering (**AAAI 2016**)
- Efficient k-Support-Norm Regularized Minimization via Fully Corrective Frank-Wolfe Method (**IJCAI 2016**)
- Efficient λ^2 Kernel Linearization via Random Feature Maps (**IEEE T-NNLS 2016**)
- Blessing of Dimensionality: Recovering Mixture Data via Dictionary Pursuit, (**IEEE T-PAMI 2016**)



Dual sparse constrained cascade regression model (IEEE T-IP 2015)

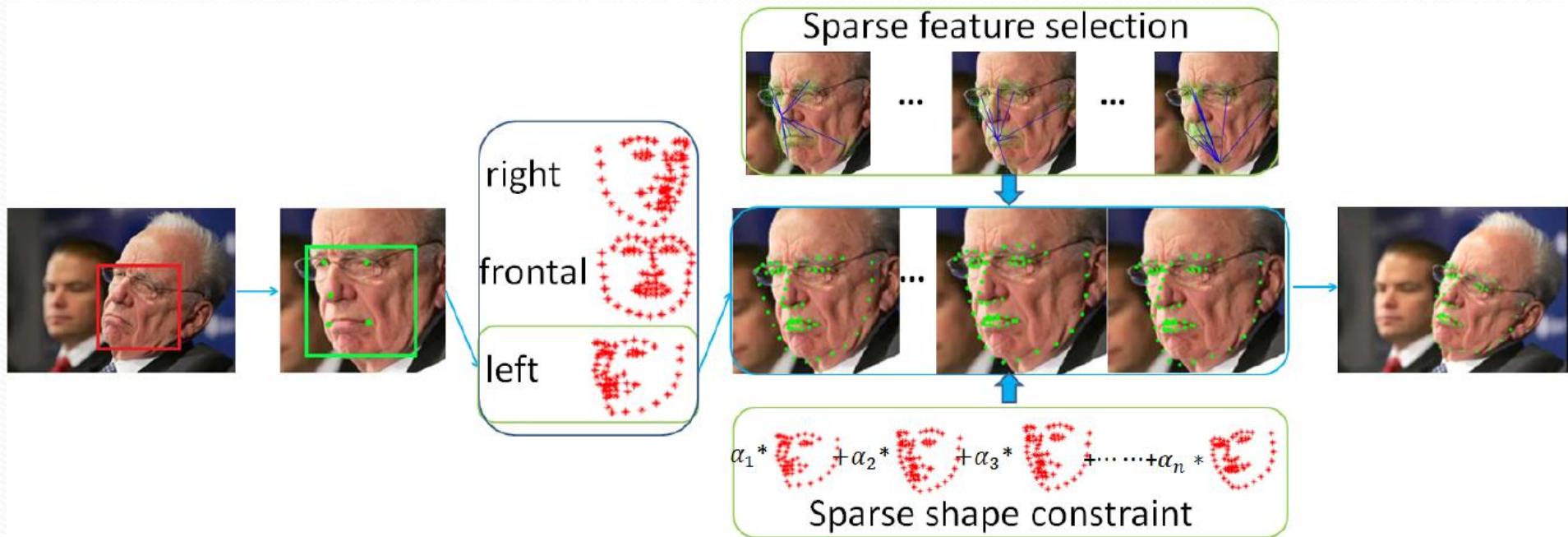


CSR:
$$\arg \min_{W_t} \sum_{i=1}^N \| (X_i^* - X_i^{t-1}) - W_t \Phi(I_i, X_i^{t-1}) \|_2^2$$

D. Piotr, W. Peter, and P. Pietro. Cascaded pose regression. *Intl. Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2010.



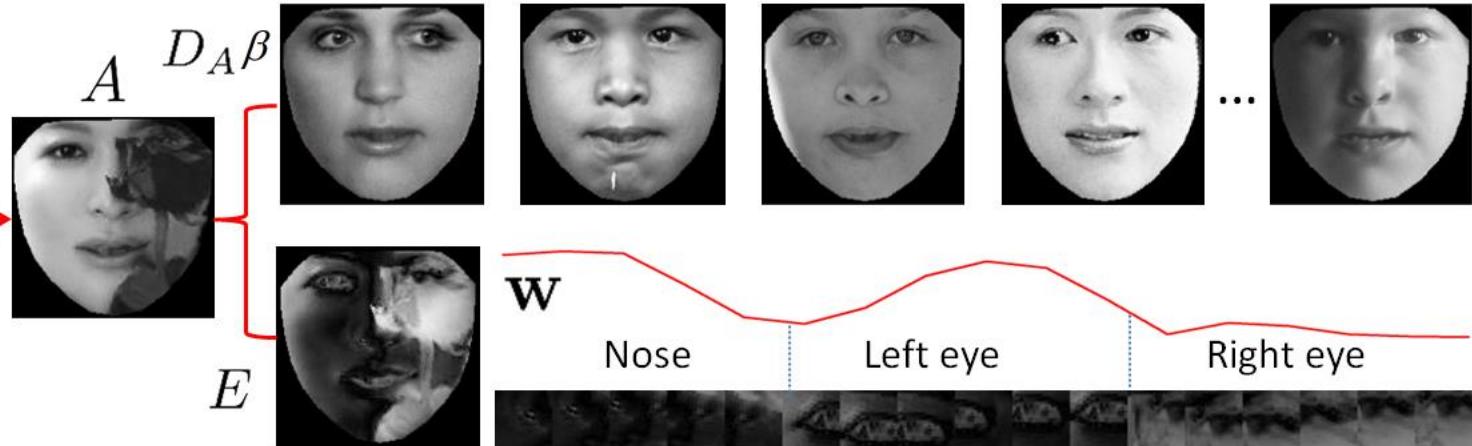
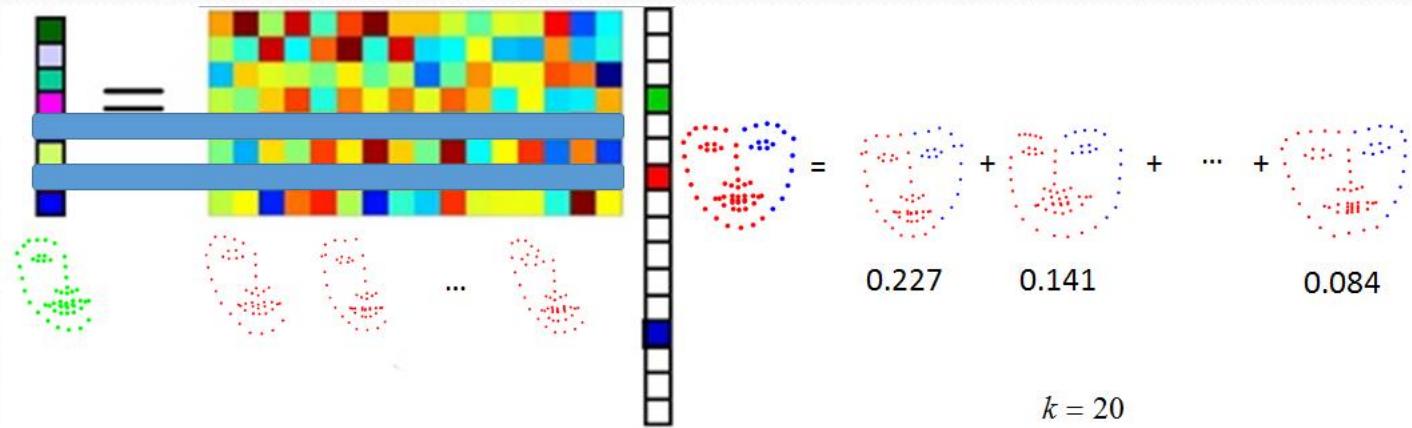
Dual sparse constrained cascade regression model (IEEE T-IP 2015)



$$\arg \min_{\alpha, \gamma, W} \|X^* - \Psi(D\alpha, \gamma) - W\Phi(I, \Psi(D\alpha, \gamma))\|_2^2 + \lambda_1 \|W\|_1 + \lambda_2 \|\alpha\|_1$$



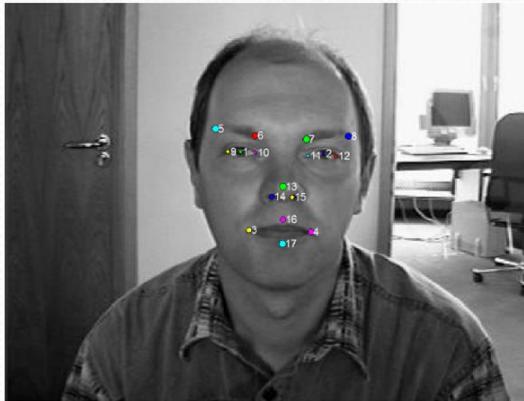
Face Alignment



Results



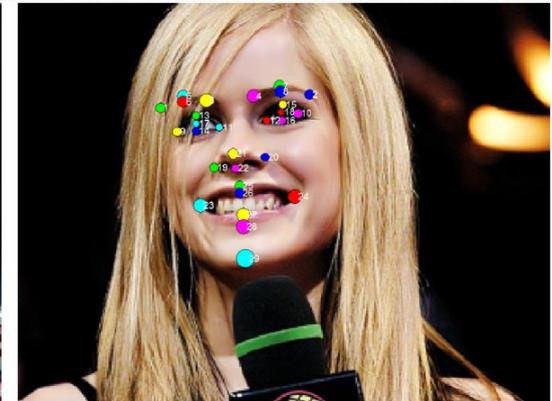
LFW



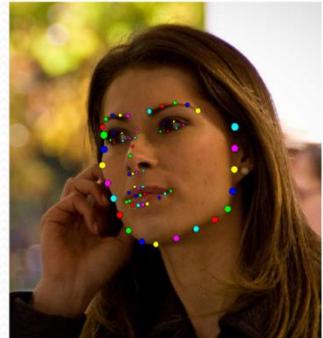
BiOID



LFPW



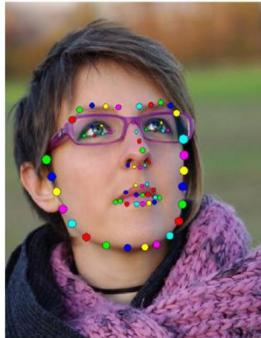
COFW



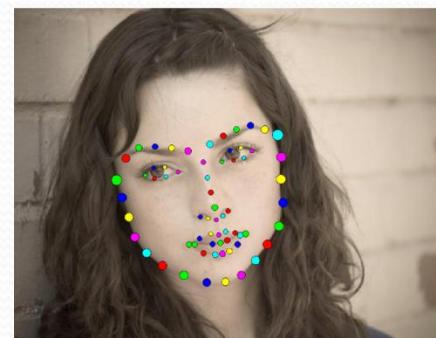
Common



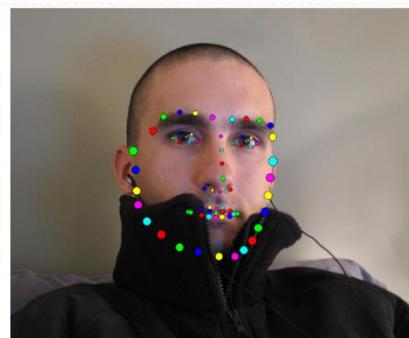
Challenge



FULL



MVFW



OCFW

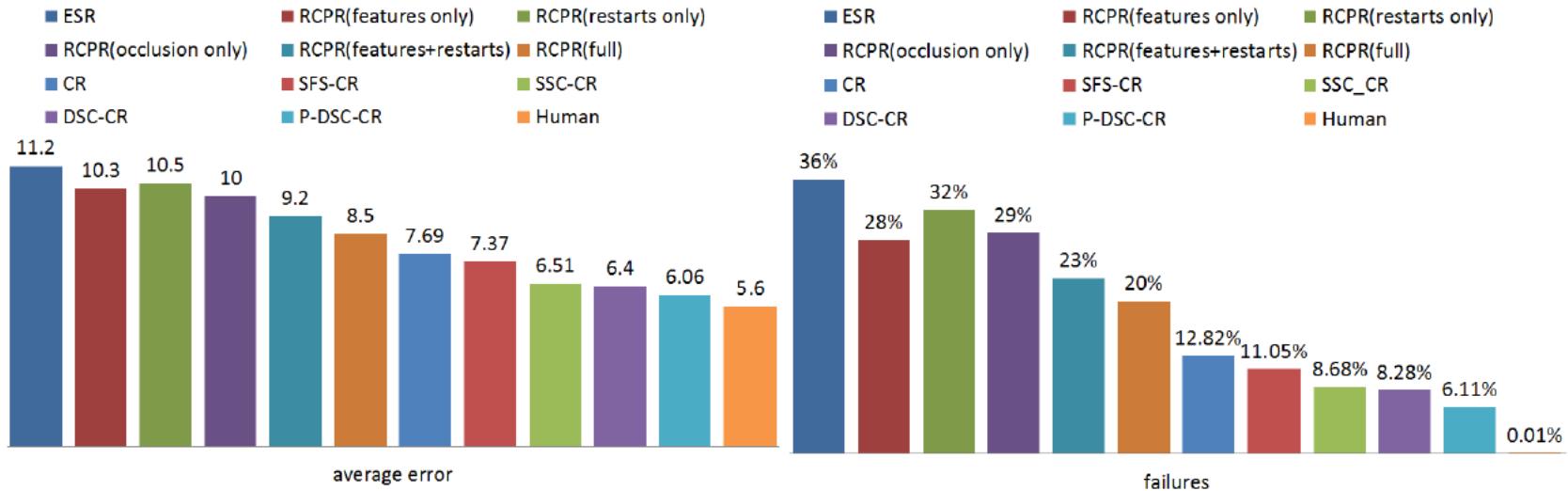


南京信息工程大学



江苏省大数据分析技术重点实验室
Jiangsu Key Laboratory of Big Data Analysis Technology

Results



Normalized mean error and the failure rate on the COFW dataset

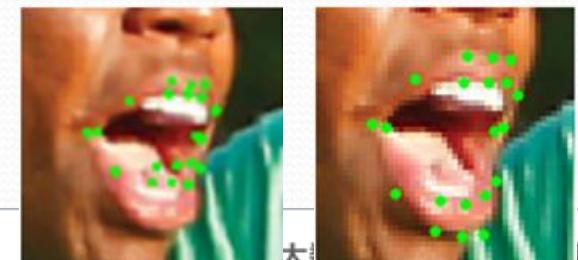
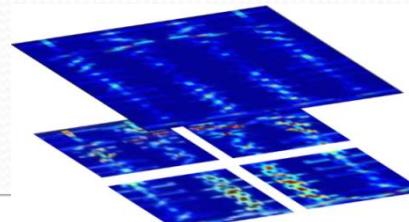
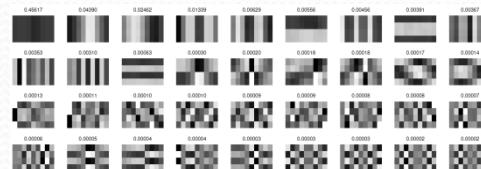
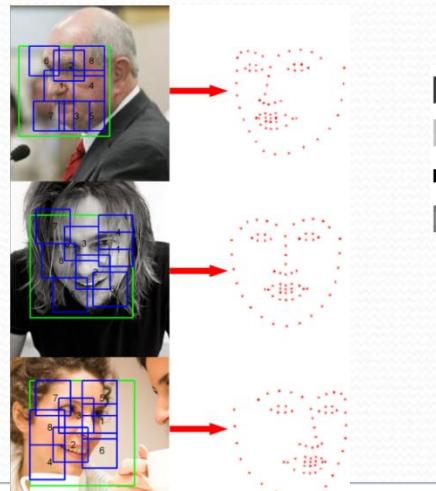
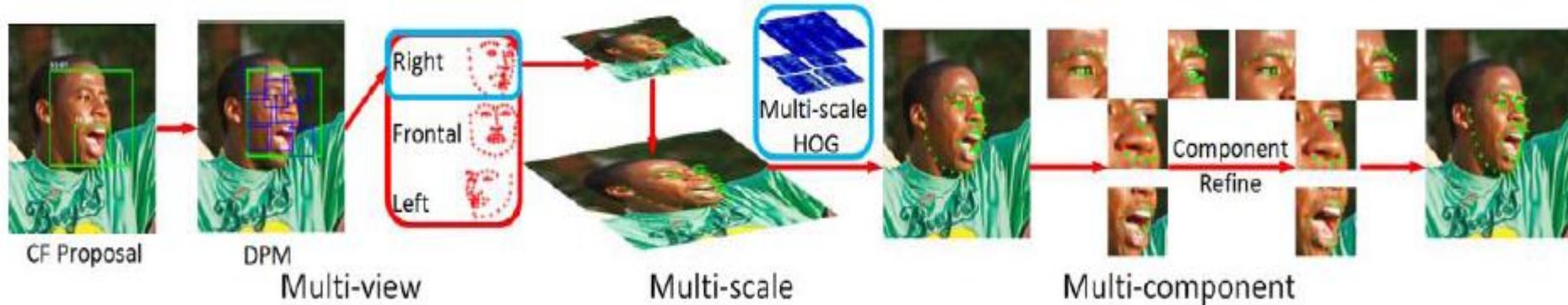
	ESR	SDM	LBF	LBF fast	TCDCN	TCDCN-Averaged	DSC-CR	P-DSC-CR
Common Subset	5.28	5.60	4.95	5.38	6.10	5.59	4.88	3.83
Challenging Subset	17.00	15.40	11.98	15.50	9.88	9.15	11.49	6.93
Fullset	7.58	7.52	6.32	7.37	6.83	6.29	6.04	4.38

MEAN ALIGNMENT ERRORS ON THE 300-W COMMON SUBSET, CHALLENGING SUBSET AND FULLSET(*0.01)



M³ CSR model (IVC 2016)

■ Multi-view, multi-scale and multi-component



南京信息工程大学



Jiangsu Key Laboratory of Big Data Analysis Technology

实验室



[home](#) » [resources](#) » [300 faces in-the-wild challenge \(300-w\), imavis 2014](#)

Datasets

- [Large Scale Fa](#)
- [300 Videos in t](#)
- [Challenge & W](#)
- [Affect "in-the-w](#)
- [Facial Express](#)
- [Analysis Challe](#)
- [300 Faces In-T](#)
- [\(300-W\), IMAV](#)
- [MAHNOB-HCI](#)
- [300 Faces In-tl](#)
- [ICCV 2013](#)
- [MAHNOB Laug](#)
- [MAHNOB MHI](#)
- [Facial point an](#)
- [MMI Facial exp](#)

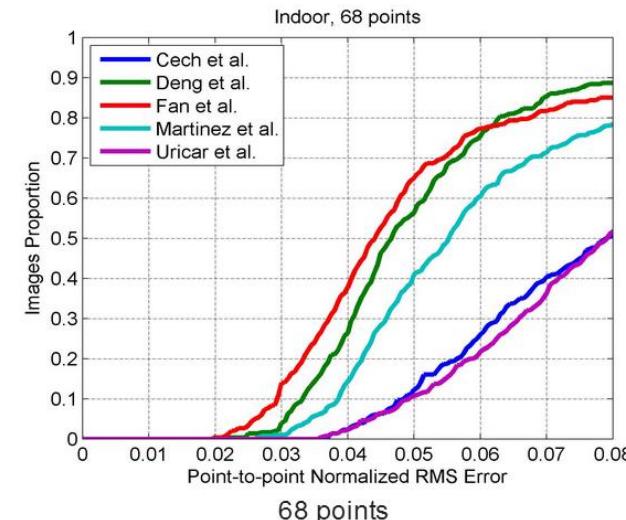
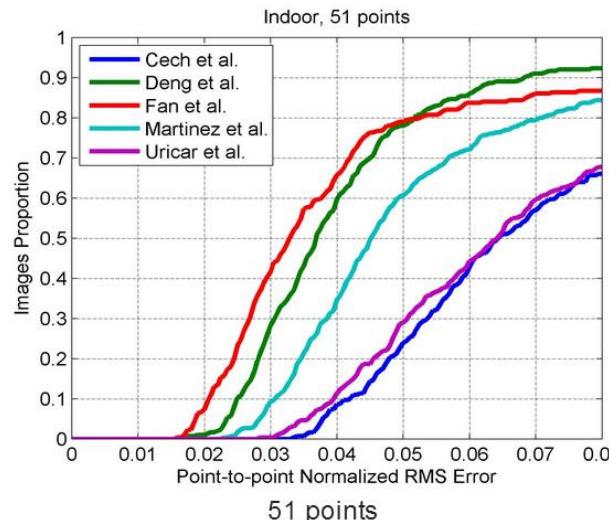
Winners

- J. Deng, Q. Liu, J. Yang, D. Tao. **M3 csr: Multi-view, multi-scale and multi-component cascade shape regression.** (Academia)
- H. Fan, E. Zhou. **Approaching human level facial landmark localization by deep learning.** (Industry)

VIS 2014

Results

Indoor



ings (secs)
12.9
0.17
1.97
1.29
2.46
5.81
42.5
12.6
3.46
4.05
—



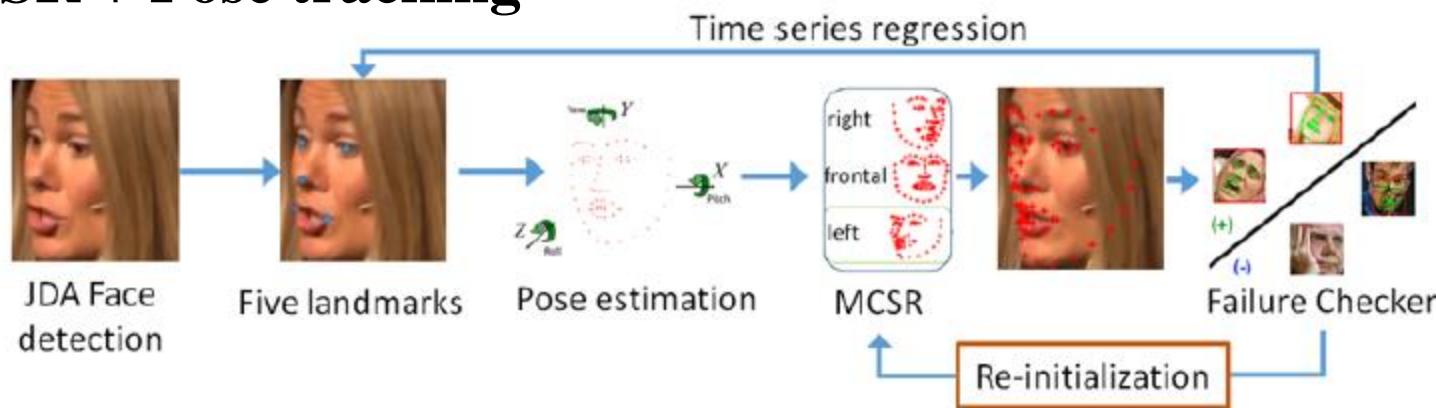
南京信息工程大学



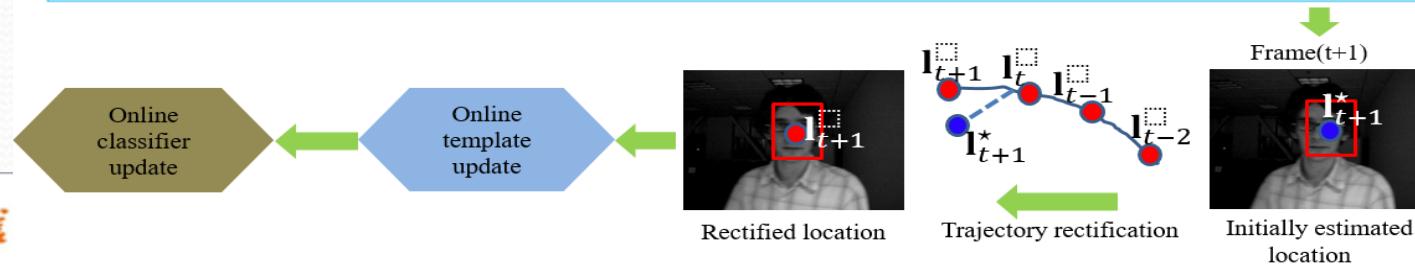
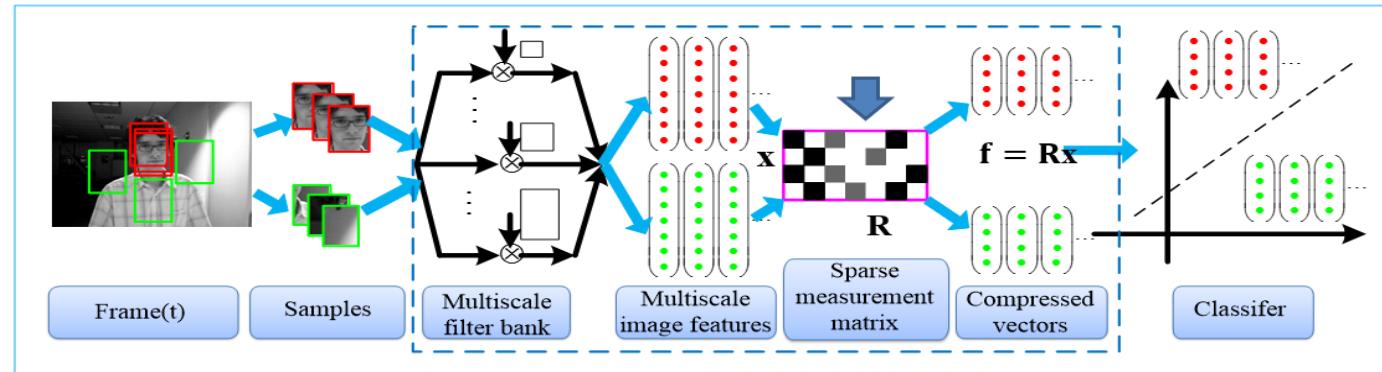
江苏省大数据分析技术重点实验室
 Jiangsu Key Laboratory of Big Data Analysis Technology

Spatio-temporal CSR (ICCVW 2015)

CSR + Pose tracking



Adaptive compressive sensing tracker (CVIU / IEEE T-CYB 2016)

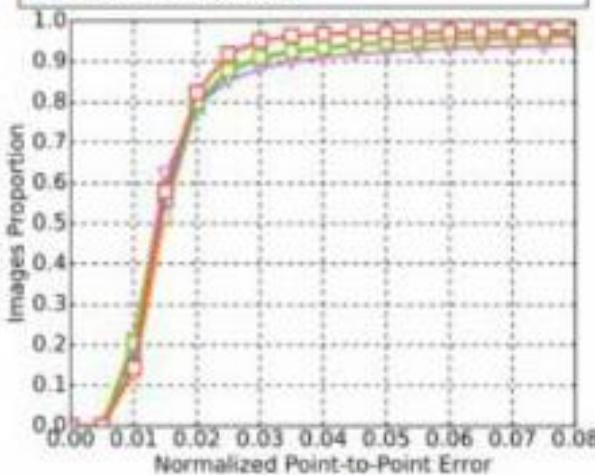
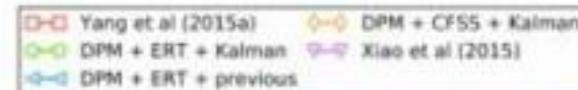




[home](#) » [resources](#) » [300 videos in the wild \(300-vw\) challenge & workshop \(iccv 2015\)](#)

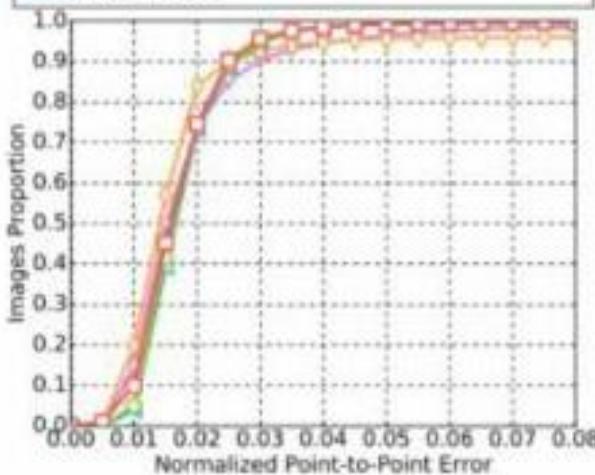
Datasets

Large Scale Facial Model

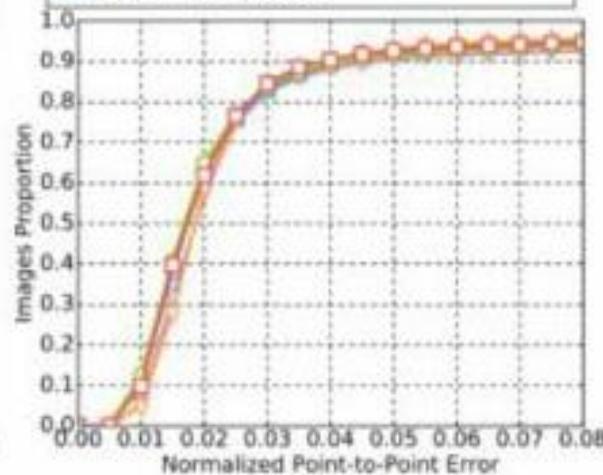


(a) Category 1

300 VIDEOS IN THE WILD (300-VW) CHALLENGE & WORKSHOP (ICCV 2015)



(b) Category 2

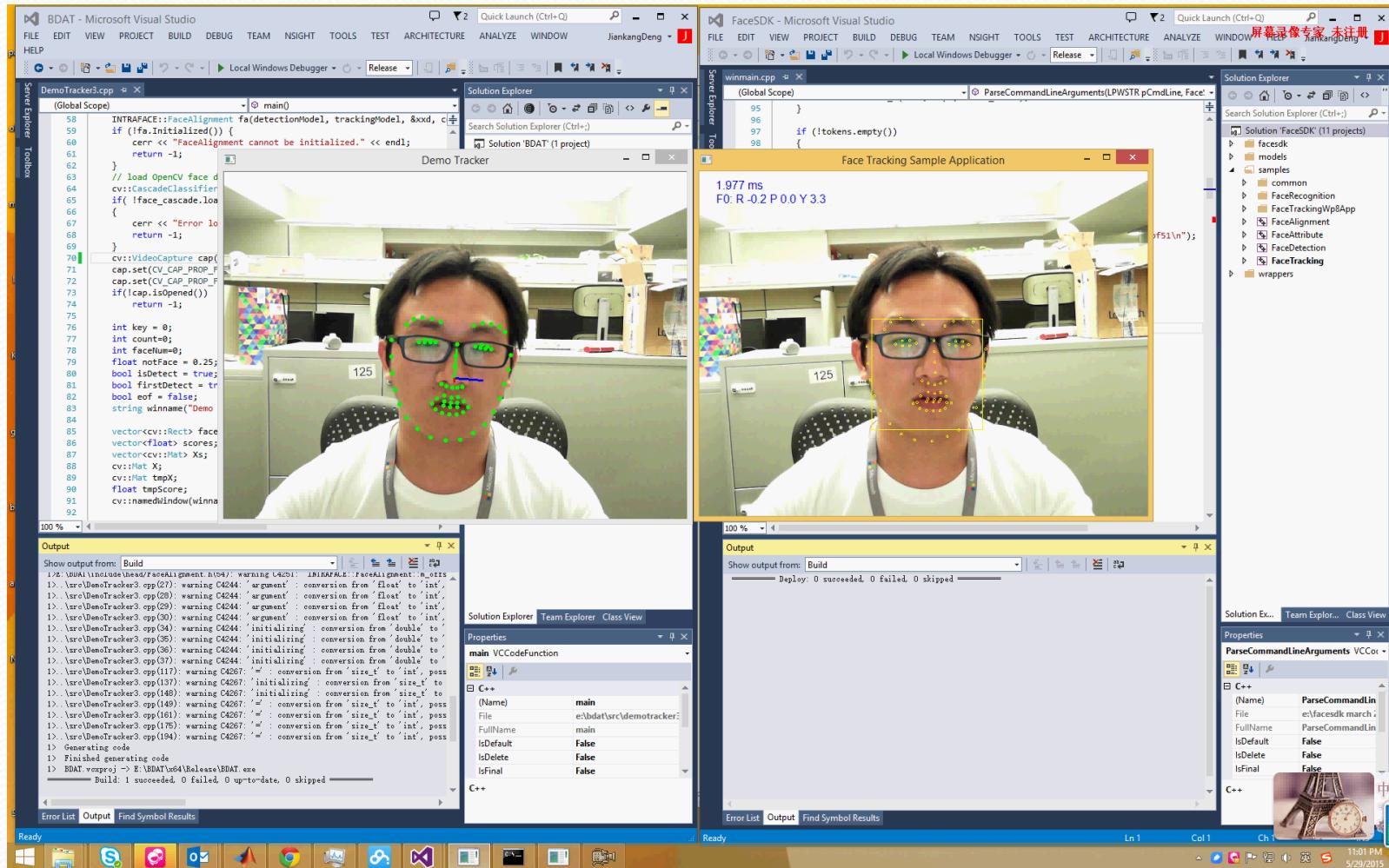


(c) Category 3

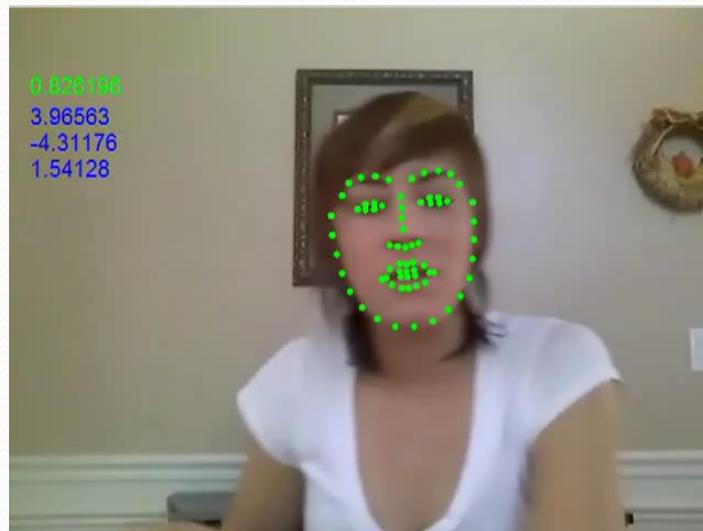
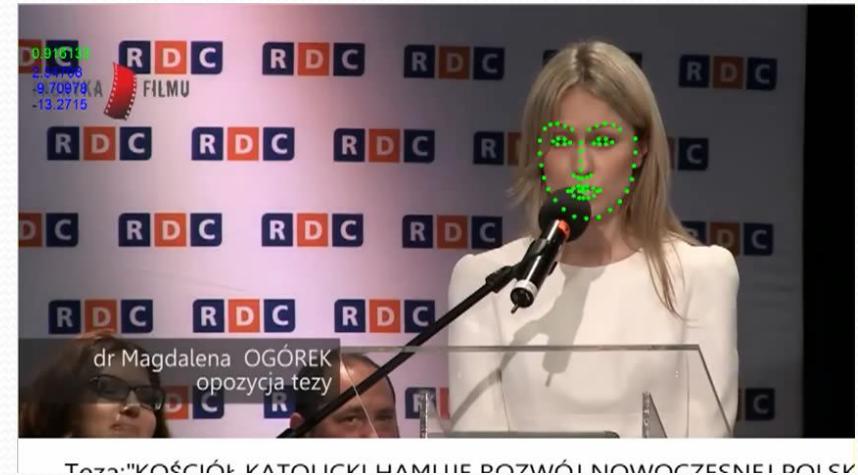
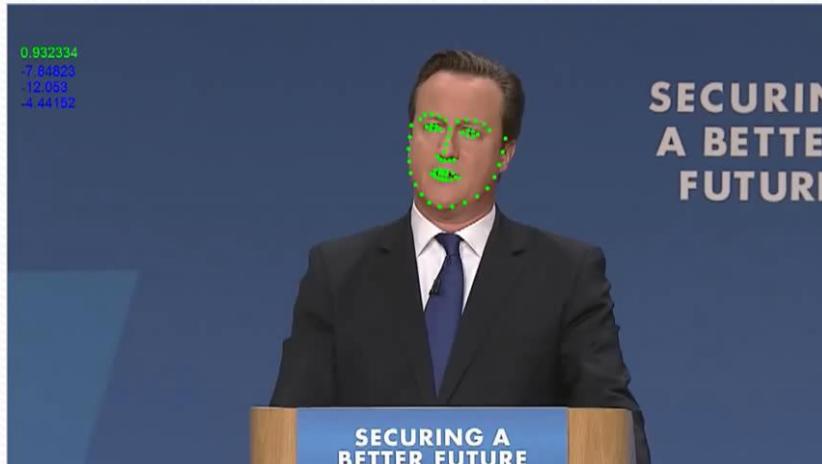
Fig. 13: Comparison between the best methods of Sections 4.3–4.7 and the participants of the 300VW challenge by Shen et al (2015). The top 5 methods are shown and are coloured red, blue, green, orange and purple, respectively. Please see Table II for a full summary.

curves are highlighted for each video category.

Video demo



Video demo



Live demo



南京信息工程大学



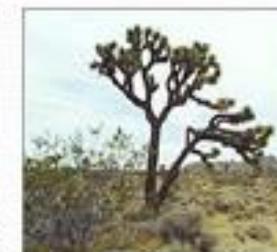
江苏省大数据分析技术重点实验室
Jiangsu Key Laboratory of Big Data Analysis Technology

Outline

- **Sparse based feature representation**
- **Hypergraph-based feature representation**
- **Deep feature representation**



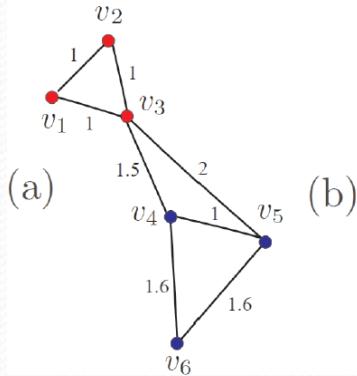
Why is hypergraph?



Six images from Caltech-101. The first three images are from the 'ferry' class; the last three images are from the 'joshua tree' class.

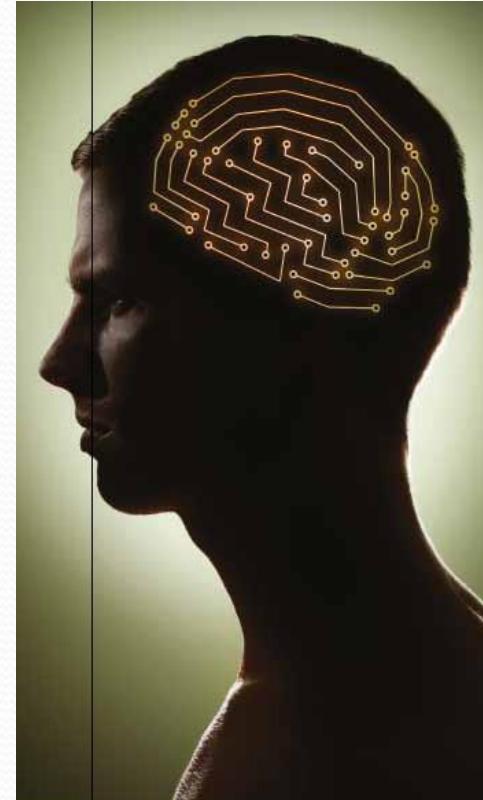
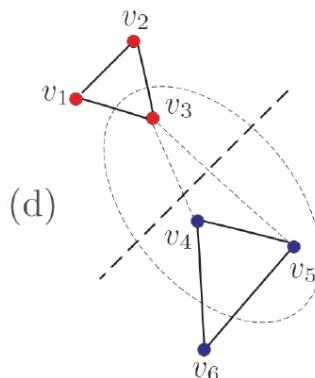
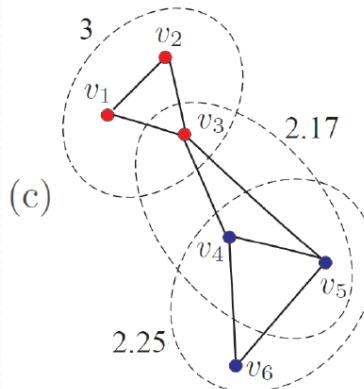
How to build the complicated relationship of multiple features?

Why is hypergraph?



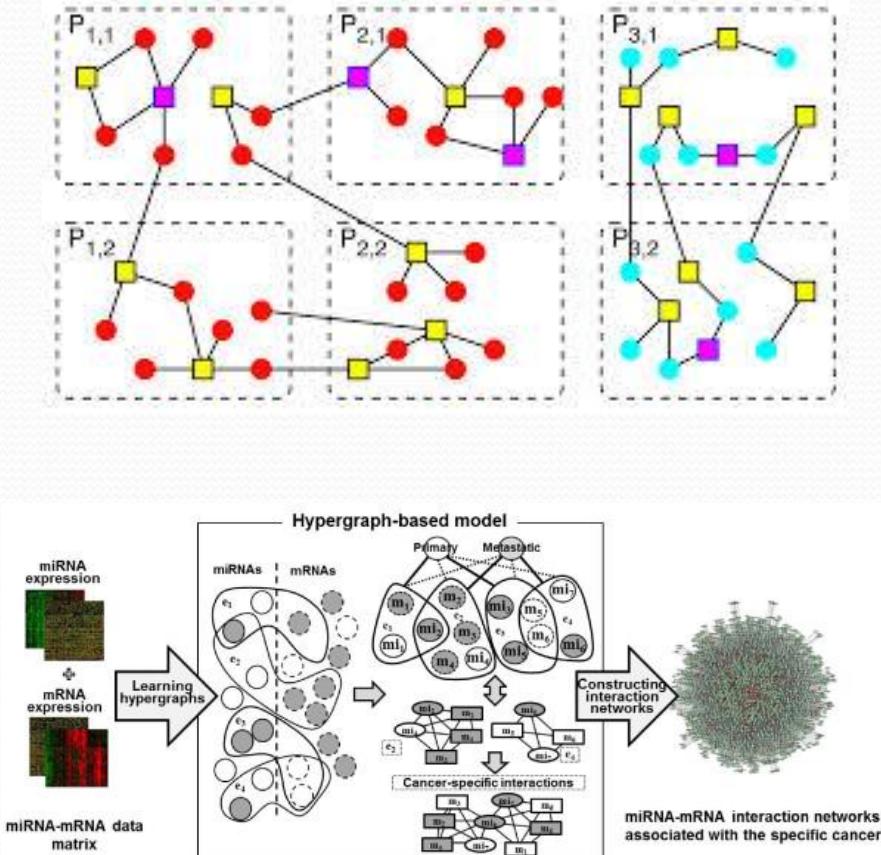
(b)

	e_1	e_2	e_3	e_4	e_5	e_6
v_1	1	1	1	0	0	0
v_2	1	1	1	0	0	0
v_3	1	1	1	1	0	0
v_4	0	0	0	1	1	1
v_5	0	0	0	1	1	1
v_6	0	0	0	0	1	1



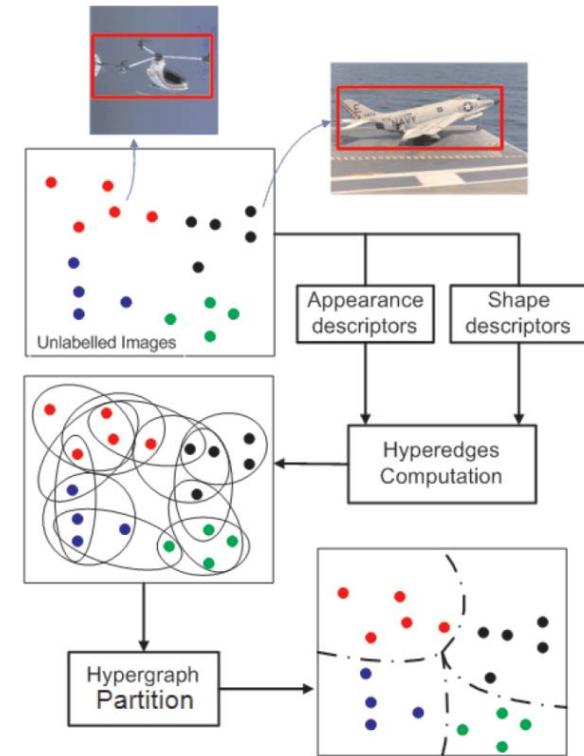
- It is not complete to represent the relations among vertices only by pairwise simple graphs.
- It may be helpful to take account of the relationship not only between two vertices, but also among three or more vertices containing local grouping information.

Why is Hypergraph?

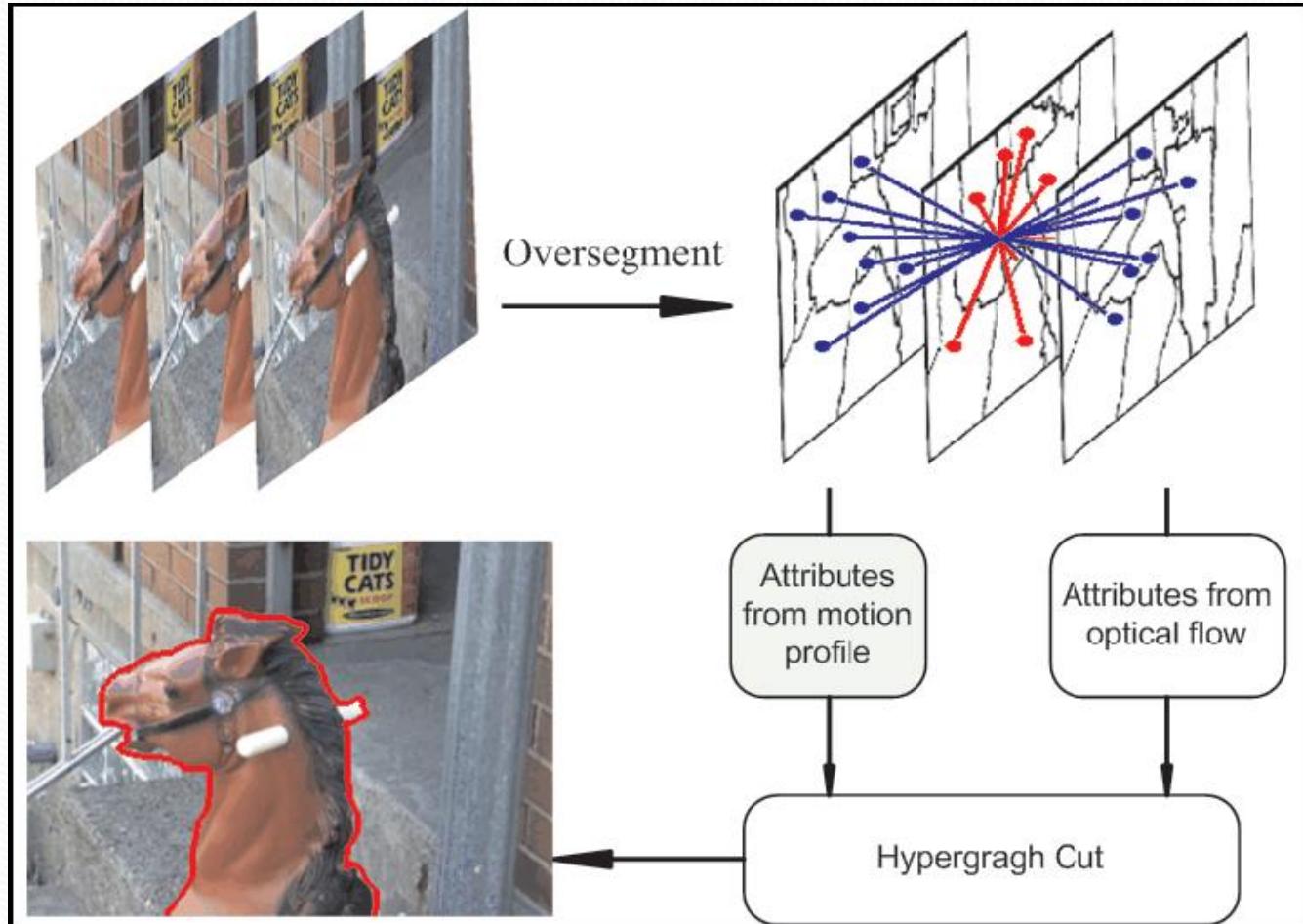


Hypergraph-based feature representation

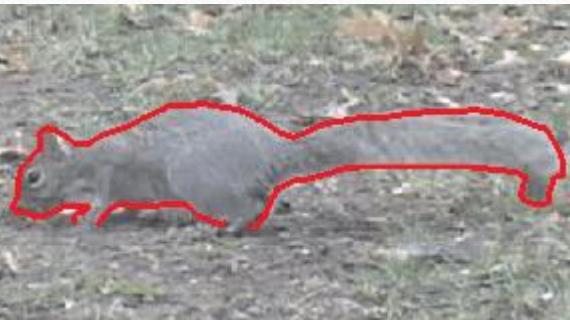
- Unsupervised hypergraph learning
 - Video objects clustering (**CVPR 2009**)
 - Image categorization (**TPAMI 2011**)
- Semi-supervised hypergraph learning
 - Content-based image retrieval (**CVPR 2010, PR 2011**)
- Sparse hypergraph learning
 - Elastic hypergraph (**TIP 2016**)
 - Application in hyperspectral image classification (TGRS submitted)



Video Object Segmentation (ICCV 2009)



Results-Squirrel



Ground Truth



Simple Graph + Optical Flow



Simple Graph + Motion Profile

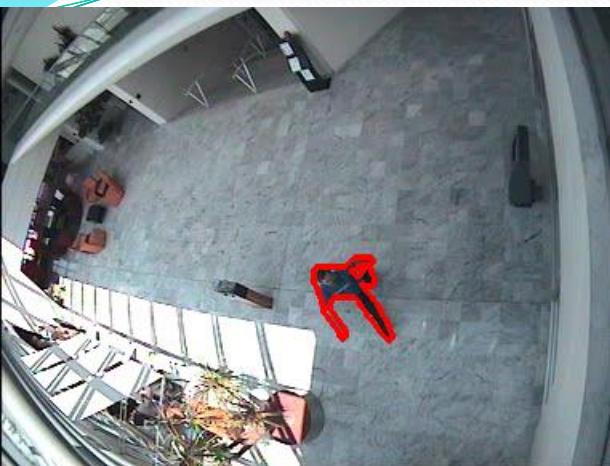


Simple Graph + Both Motion Cues



Hypergraph Cut

Results-Walking with Rotation



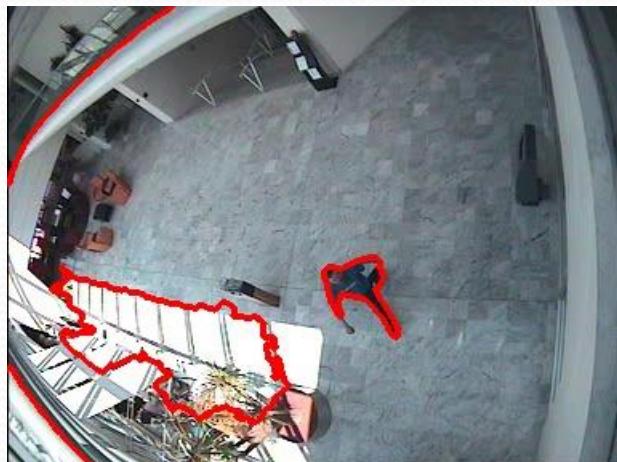
Ground Truth



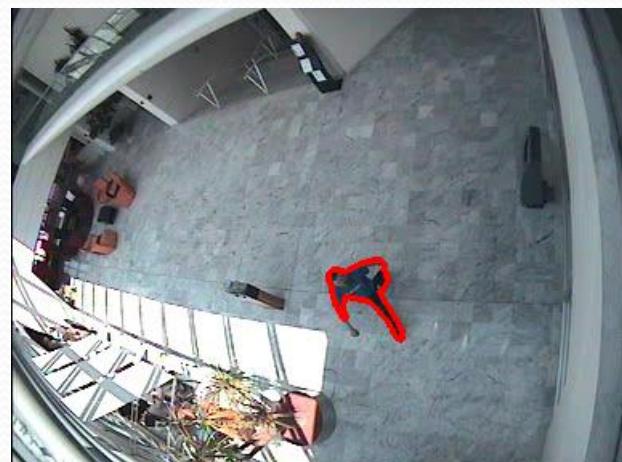
Simple Graph + Optical Flow



Simple Graph + Motion Profile

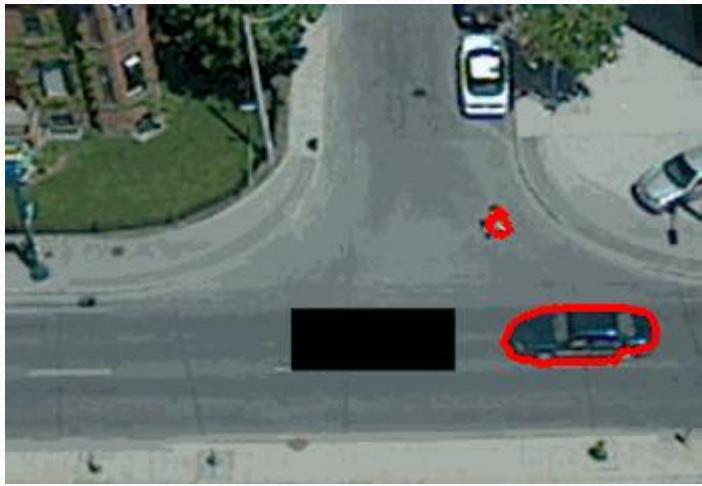


Simple Graph + Both Motion Cues



Hypergraph Cut

Videos

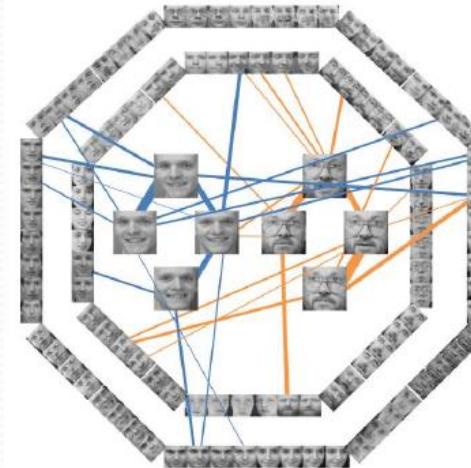
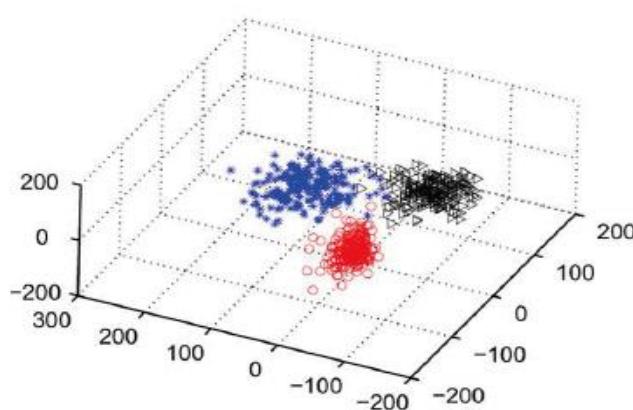


Elastic Net Hypergraph Learning (IEEE T-IP 2016)

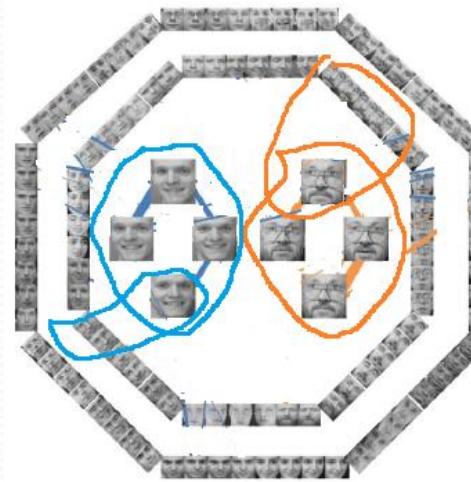
Robust Elastic Net Representation

$$\begin{aligned} & \min_{Z,S} \|Z\|_1 + \lambda \|Z\|_F^2 + \gamma \|S\|_{2,1} \\ & \text{s.t. } X = XZ + S, \text{diag}(Z) = 0 \end{aligned}$$

Hypergraph Learning



KNN-Graph



Elastic Net
Hypergraph

Elastic Net Hypergraph Learning (IEEE T-IP 2016)

Dataset	G-graph	LE-graph	l_1-graph	KNN-HG	SCHG	l_1-Hypergraph	ENHG
Extended Yale B (10%)	66.49	70.79	53.87	71.80	77.68	82.15	88.59
Extended Yale B (20%)	65.34	69.97	54.46	75.54	81.80	83.48	90.87
Extended Yale B (30%)	33.72	71.85	53.90	77.67	82.84	85.36	93.94
Extended Yale B (40%)	66.28	71.34	56.61	80.59	83.55	86.90	94.34
Extended Yale B (50%)	66.90	71.60	57.75	80.80	84.48	87.08	94.97
Extended Yale B (60%)	67.52	71.48	58.48	81.79	89.46	90.42	95.28
PIE (10%)	65.72	67.75	78.29	68.74	79.35	80.24	88.31
PIE (20%)	66.94	69.58	82.82	70.18	84.74	84.55	94.94
PIE (30%)	69.89	73.48	87.94	74.39	88.78	89.29	96.55
PIE (40%)	71.54	76.38	90.99	76.14	90.33	91.75	97.33
PIE (50%)	73.04	78.35	93.39	78.76	92.66	93.71	97.53
PIE (60%)	74.91	80.44	95.00	79.95	94.12	94.87	98.43
USPS (10%)	96.87	96.79	88.33	96.51	97.08	97.20	97.36
USPS (20%)	97.78	97.90	91.11	98.17	98.12	98.29	98.39
USPS (30%)	98.45	98.47	93.08	98.78	98.87	98.85	98.91
USPS (40%)	98.80	98.82	95.96	99.08	99.08	99.10	99.08
USPS (50%)	99.18	99.14	97.31	99.39	99.41	99.39	99.40
USPS (60%)	99.35	99.28	98.86	99.51	99.50	99.52	99.53

Outline

- Sparse based feature representation
- Hypergraph-based feature representation
- Deep feature representation



Deep Learning



Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton* and R. R. Salakhutdinov

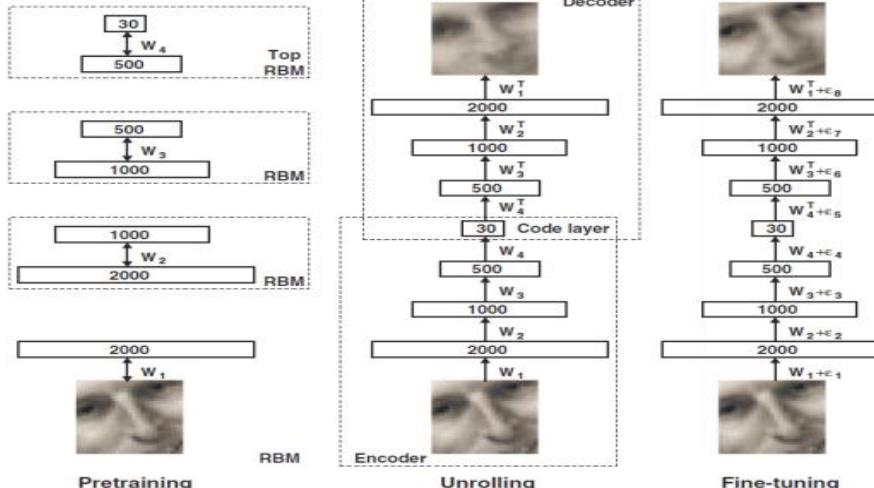
2006

High-dimensional data can be converted to low-dimensional codes by training a neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such "autoencoder" networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

Dimensionality reduction facilitates the classification, visualization, communication, and storage of high-dimensional data. A simple and widely used method is principal components analysis (PCA), which

finds the directions of greatest variance in the data set and represents each data point by its coordinates along each of these directions. We describe a nonlinear generalization of PCA that uses an adaptive, multilayer "encoder" network

2006 VOL 313 SCIENCE www.sciencemag.org



Breakthrough



2011年

Speech recognition

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	24	30.4	36.2
(Sentence error rates)			
Google Voice Input	5,870	12.3	
Youtube	1,400	47.6	52.3

2012年

Image classification

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

No. 1 in 10 breakthrough tech 2013 selected by MIT tech review

MIT Technology Review

10 BREAKTHROUGH TECHNOLOGIES 2013

Introduction The 10 Technologies Past Years

Deep Learning	Temporary Social Media	Prenatal DNA Sequencing	Additive Manufacturing	Baxter: The Blue-Collar Robot
With massive amounts of computational power, machines can now recognize objects and translate speech in real time. Artificial intelligence is finally getting smart.	Messages that quickly self-destruct could enhance the privacy of online communications and make people freer to be spontaneous.	Reading the DNA of fetuses will be the next frontier of the genomic revolution. But do you really want to know about the genetic problems or musical aptitude of your unborn child?	Skeptical about 3-D printing? GE, the world's largest manufacturer, is on the verge of using the technology to make jet parts.	Rodney Brooks's newest creation is easy to interact with, but the complex innovations behind the robot show just how hard it is to get along with people.
Memory Implants	Smart Watches	Ultra-Efficient Solar Power	Big Data from Cheap Phones	Supergrids
A maverick neuroscientist believes he has deciphered the code by which memory forms long-term memories. Next: testing a prosthetic implant for people suffering from long-term memory loss.	The designers of the Pebble watch realized that a mobile phone is more useful if you don't have to take it out of your pocket.	Doubling the efficiency of a solar panel could completely change the economics of renewable energy. Nanotechnology just might make it possible.	Collecting and analyzing information from simple cell phones could provide surprising insights into how people move about and behave – and even help us understand the spread of diseases.	A new high-power circuit breaker could finally make highly efficient DC power grids practical.

REVIEW

doi:10.1038/nature14539

Deep learning

Yann LeCun^{1,2}, Yoshua Bengio³ & Geoffrey Hinton^{4,5}

Deep learning allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction. These methods have dramatically improved the state-of-the-art in speech recognition, visual object recognition, object detection and many other domains such as drug discovery and genomics. Deep learning discovers intricate structure in large data sets by using the backpropagation algorithm to indicate how a machine should change its internal parameters that are used to compute the representation in each layer from the representation in the previous layer. Deep convolutional nets have brought about breakthroughs in processing images, video, speech and audio, whereas recurrent nets have shone light on sequential data such as text and speech.

2015年5月Nature杂志以综述的形式对深度学习进行了总结和评价，指出深度学习最大的优点是能自动学习和抽象数据特征。

机器学习是研究如何使计算机通过输入输出模式识别新输入的子领域。现代机器学习的一个重要方面是深度学习，它允许计算机自动地从大量数据中学习。深度学习在许多应用中取得了显著的进展，包括语音识别、图像识别、自然语言处理等。深度学习的成功归功于其强大的表示学习能力，能够发现数据中的高级结构。然而，深度学习仍然面临着一些挑战，如模型的解释性和泛化能力。未来的研究方向可能包括改进模型的解释性、提高泛化能力和探索新的深度学习架构。

IMAGENET Large Scale Visual Recognition Challenge 2015 (ILSVRC2015)

Task 1b: Object detection with additional training data

[News](#)[Home](#)

Ordered by number of categories won

[Sources](#)[Registration](#)[FAQ](#)[Citation](#)[Contact](#)**News**

- Feature extraction from images
- Jan 2015: Amax remove threshold compared to entry1
- Jan 2015: CUIimage Combined models with region proposals of cascaded RPN, edgebox and selective search
- Dec 2014: MIL-UT ensemble of 4 random forests
- Nov 2014: Amax Cascade region proposal
- Oct 2014: Ensemble of 4 random forests learned separately
- Oct 2014: Amax
- August 15, 2015: Registration deadline
- August 13, 2015: Competition deadline
- June 12, 2015: Tentative competition date
- June 2, 2015: Additional competition date
- May 19, 2015: Announcement of competition date
- December 17, 2015: Status report due
- September 19, 2014: Training data released in Chapel Hill.

Team name	Entry description	Description of outside data used	Number of object categories won	mean AP
Amax	remove threshold compared to entry1	pre-trained model from classification task; add training examples for class number <1000	165	0.57848
CUIimage	Combined models with region proposals of cascaded RPN, edgebox and selective search	3000-class classification images from ImageNet are used to pre-train CNN	30	0.522833

Task 3b: Object detection from video with additional training data

Ordered by number of categories won

Team name	Entry description	Description of outside data used	Number of object categories won	mean AP
Amax	only half of the videos are tracked due to deadline limits, others are only detected by Faster RCNN (VGG16) without temporal smooth.	---	18	0.730746
CUVideo	Outside training data (ImageNet 3000-class data) to pre-train the detection model, mAP 77.0 on validation data	ImageNet 3000-class data to pre-train the model	11	0.696607
Tramps-Soushen	Models combine with mAP constraint			
Tramps-Soushen	Combine several mode			
BAD	VID2015_trace_merge			
BAD	combined_VIDtrainval			
BAD	VID2015_merge_test_t			
BAD	combined_test_DET_th			
BAD	VID2015_VID_test_thr			

Results announced.

Task 2b: Classification+localization with additional training data

Ordered by classification error

Team name	Entry description	Description of outside data used	Classification error	Localization error
Amax	Validate the classification model we used in DET entry1	share proposal procedure with DET for convinence	0.04354	0.14574
Tramps-Soushen	extra annotations collected by ourselves	extra annotations collected by ourselves	0.04581	0.122285
CUIimage	Average multiple models. Validation accuracy is 79.78%.	3000-class classification images from ImageNet are used to pre-train CNN	0.05858	0.198272

History[2014](#), [2013](#), [2012](#), [2011](#), [2010](#)



IMAGENET Large Scale Visual Recognition Challenge 2016 (ILSVRC2016)

Object detection from video (VID)^[top]

[News](#) [History](#)

Task 3a: Object detection from video with provided training data

News

Ordered by number of categories won

- September 18, 2016 5pm PST.
- May 31, 2016: CUVi
- May 26, 2016: NUIST
- May 3, 2016: S

Team name	Entry description	Number of object categories won	mean AP
NUIST	cascaded region regression + tracking	10	0.808292
NUIST	cascaded region regression + tracking	10	0.803154
CUVi	1-model ensemble with Multi-Context Suppression and Motion-Guided		

Task 3b: Object detection from video with additional training data

History

[2015](#), [2014](#), [2013](#), [2012](#)

Tentative Time

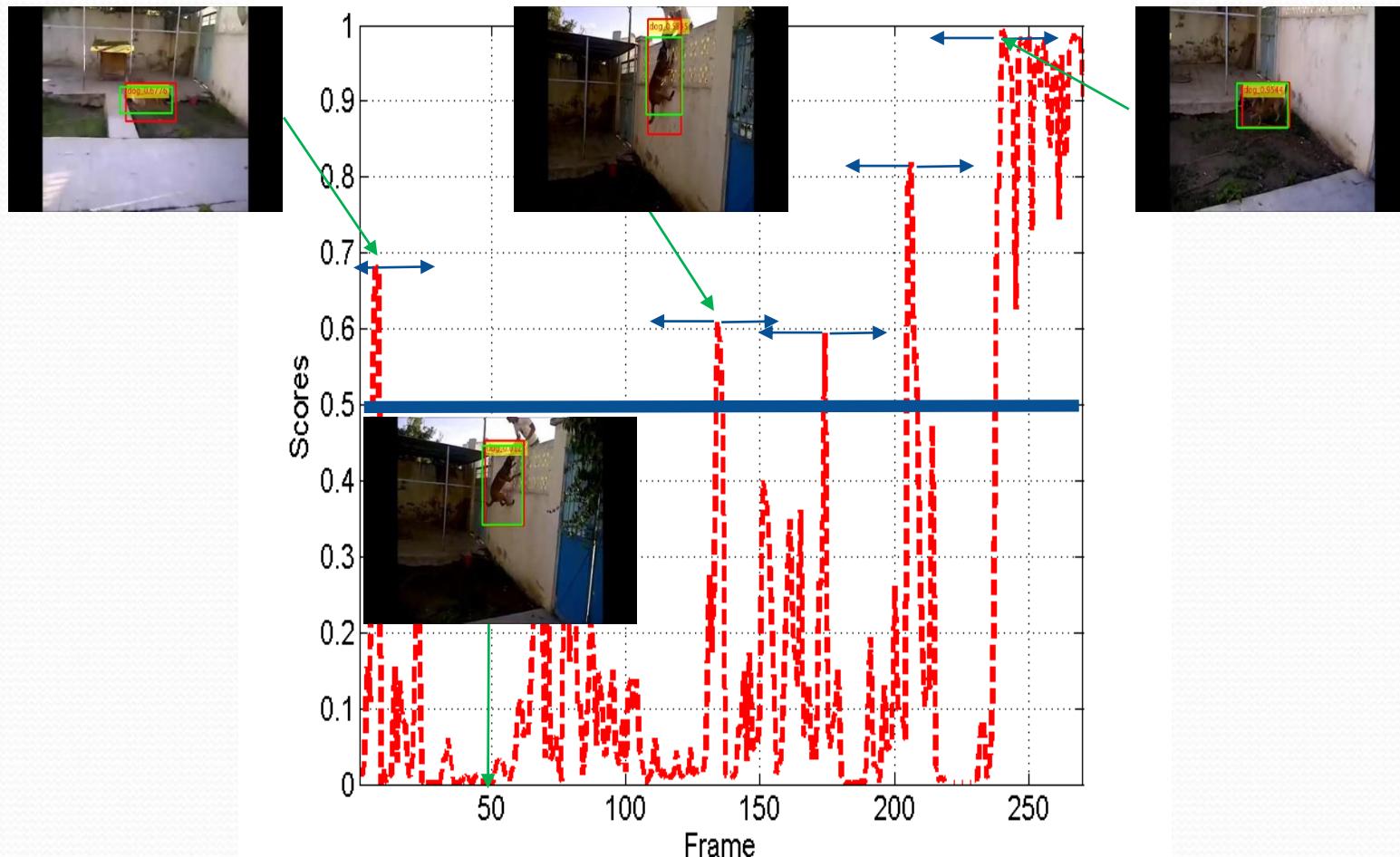
- May 31, 2016:
- September 9, 2016:

Team name	Entry description	Description of outside data used	Number of object categories won	mean AP
NUIST	cascaded region regression + tracking	proposal network is finetuned from COCO	17	0.79593
NUIST	cascaded region regression + tracking	proposal network is finetuned from COCO	5	0.781144
Trimpes-Sous		Extra data from ImageNet dataset/out of the		

Task 3d: Object detection/tracking from video with additional training data

Team name	Entry description	Description of outside data used	mean AP
NUIST	cascaded region regression + tracking	proposal network is finetuned from COCO	0.583898
ITLab-	An ensemble for detection,	pre-trained model from COCO detection, extra data collected	0.100863

Object detection from Video

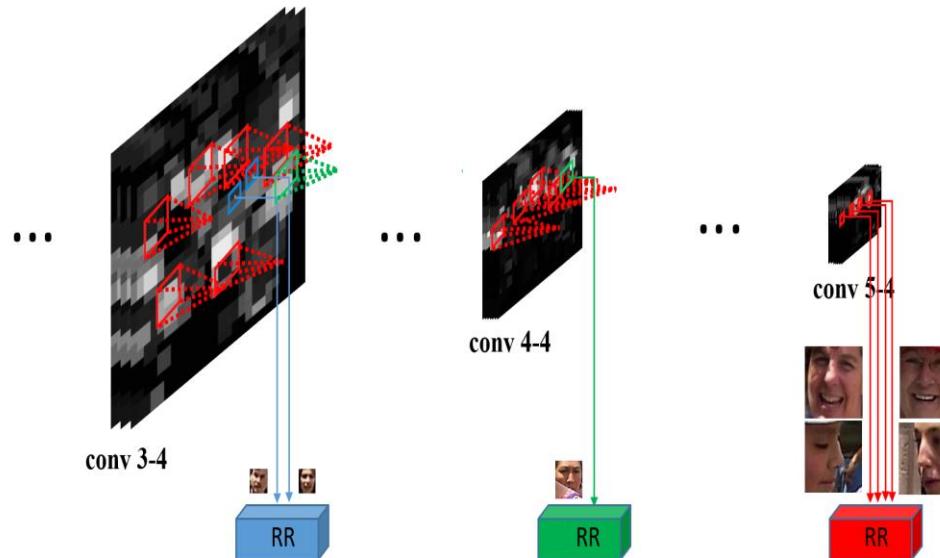


Object detection on each frame

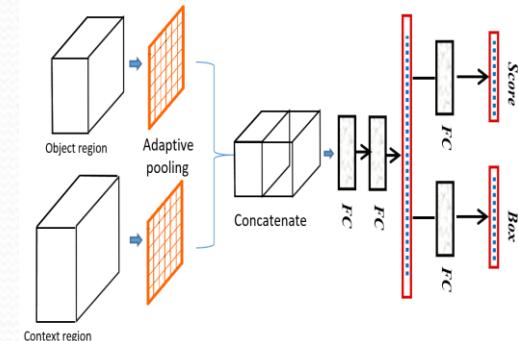
Tracking from the high score frame (temporal smooth)

Class-wise box regression and NMS on each frame

Cascade Region Regression



**Multi-layer Conv Feature
(region size specific)**



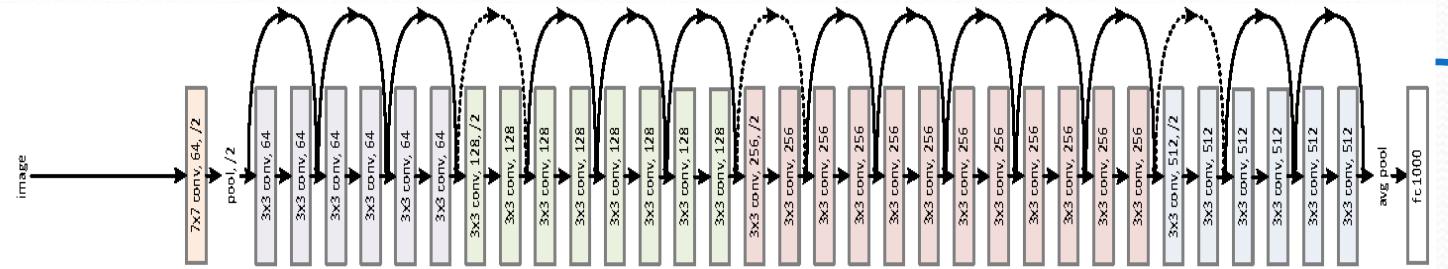
**Multi-scale Conv
Feature
(object + around
context)**

Cascade region regression根据region的大小选择不同的region regressor对bounding box进行调整，较大的region使用后面的feature map，较小的feature map使用前面的feature map。

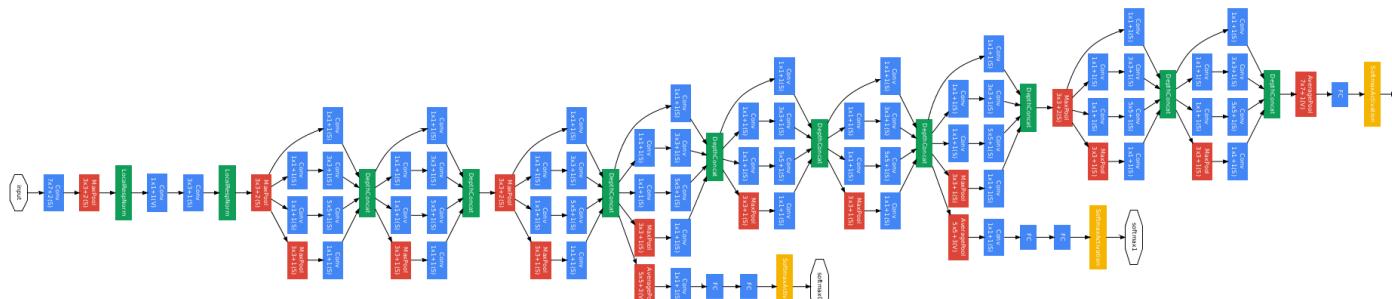
Model Ensemble

Res-net

34-layer residual



Google-net

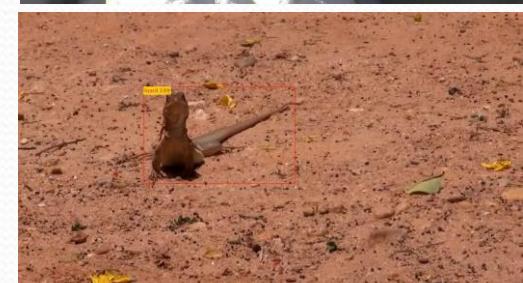
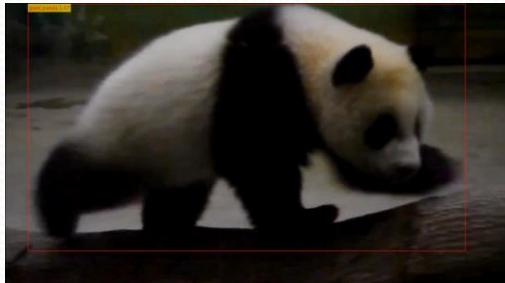


**Model ensemble
is always
effective.**

Demo Video

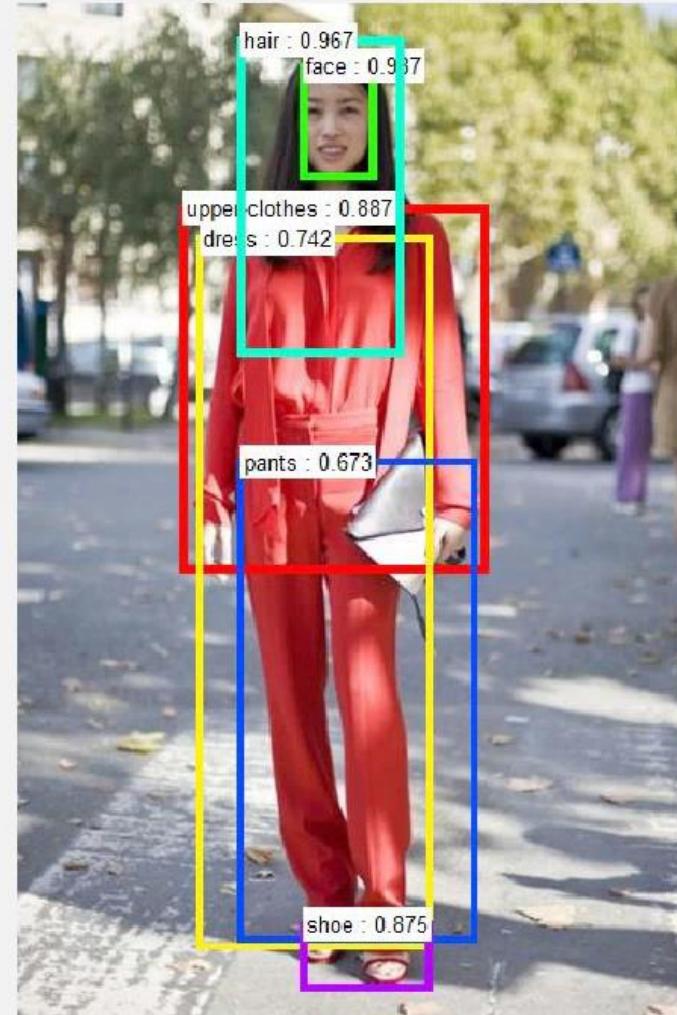


Demo Video

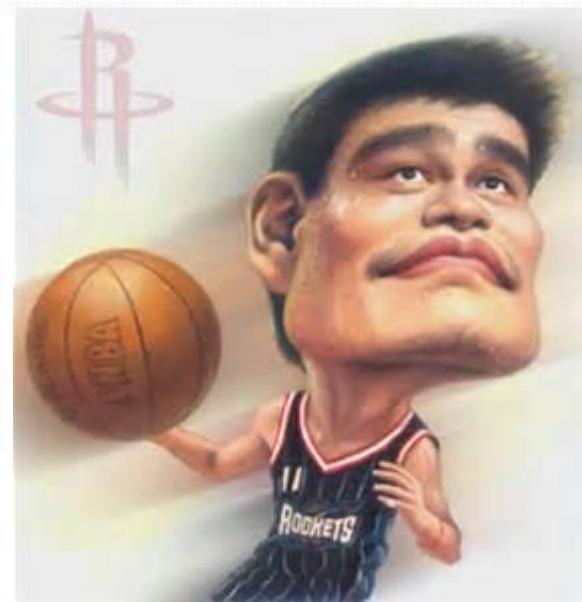


load

detect



Does Cartoonist use deep features?



Thank You

Qingshan Liu

Email: qsliu@nuist.edu.cn

Cell: 13585199482