

Object recognition (part 2)

CSE P 576

Larry Zitnick (larryz@microsoft.com)

Convolutional Nets

Yann LeCun

The Courant Institute of Mathematical Sciences

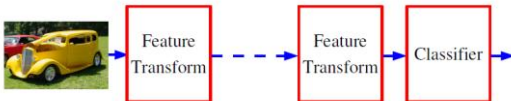
New York University

<http://yann.lecun.com>

Yann LeCun

New York University

Good Internal Representations are Hierarchical



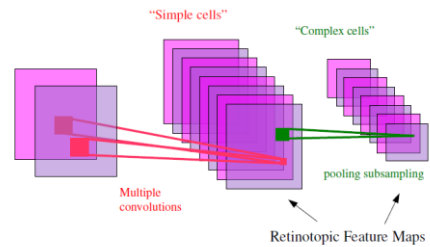
- Low-level features - mid-level features - high-level features - categories
- Representations are increasingly abstract, global, and invariant.
- In Vision: part-whole hierarchy
 - ▶ Pixels->Edges->Textons->Parts->Objects->Scenes
- In Language: hierarchy in syntax and semantics
 - ▶ Words->Parts of Speech->Sentences->Text
 - ▶ Objects,Actions,Attributes...-> Phrases -> Statements -> Stories

Yann LeCun

New York University

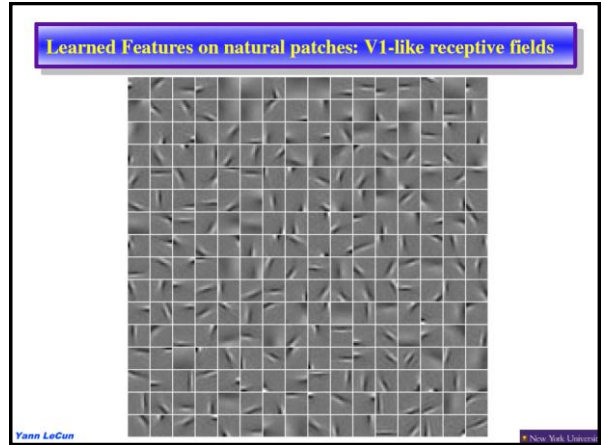
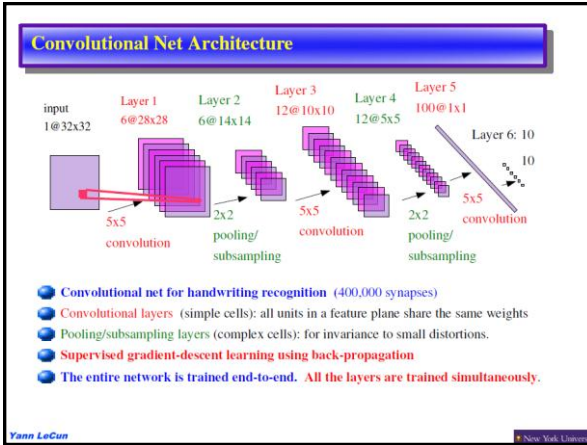
An Old Idea for Image Representation with Distortion Invariance

- [Hubel & Wiesel 1962]:
 - ▶ simple cells detect local features
 - ▶ complex cells "pool" the outputs of simple cells within a retinotopic neighborhood.



Yann LeCun

New York University



MNIST Handwritten Digit Dataset

Handwritten Digit Dataset MNIST: 60,000 training samples, 10,000 test samples


Yann LeCun New York University

Some Results on MNIST (from raw images: no preprocessing)

CLASSIFIER	DEFORMATION	ERROR	Reference
Knowledge-free methods (a fixed permutation of the pixels would make no difference)			
2-layer NN, 800 HU, CE		1.60	Simard et al., ICDAR 2003
3-layer NN, 500+300 HU, CE, reg		1.53	Hinton, in press, 2005
SVM, Gaussian Kernel		1.40	Cortes 92 + Many others
Convolutional nets			
Convolutional net LeNet-5,		0.80	Ranzato et al., NIPS 2006
Convolutional net LeNet-6,		0.70	Ranzato et al., NIPS 2006
Training set augmented with Affine Distortions			
2-layer NN, 800 HU, CE	Affine	1.10	Simard et al., ICDAR 2003
Virtual SVM deg-9 poly	Affine	0.80	Scholkopf
Convolutional net, CE	Affine	0.60	Simard et al., ICDAR 2003
Training set augmented with Elastic Distortions			
2-layer NN, 800 HU, CE	Elastic	0.70	Simard et al., ICDAR 2003
Convolutional net, CE	Elastic	0.40	Simard et al., ICDAR 2003

Note: some groups have obtained good results with various amounts of preprocessing such as deskewing (e.g. 0.56% using an SVM with smart kernels [deCoste and Schoelkopf]) hand-designed feature representations (e.g. 0.63% with "shape context" and nearest neighbor [Belongie])

Yann LeCun New York University



Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories

To appear in CVPR 2006

Svetlana Lazebnik (slazebni@uiuc.edu)
Beckman Institute, University of Illinois at Urbana-Champaign

Cordelia Schmid (cordelia.schmid@inrialpes.fr)
INRIA Rhône-Alpes, France


Jean Ponce (ponce@di.ens.fr)
Ecole Normale Supérieure, France

http://www-cvr.ai.uiuc.edu/ponce_grp

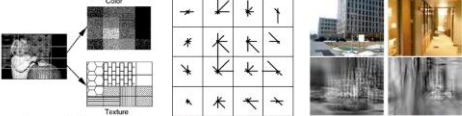
1

Overview

- A "pre-attentive" approach: recognize the scene as a whole without examining its constituent objects Biederman (1988), Thorpe et al. (1996), Fei-Fei et al. (2002), Renninger & Malik (2004)
- Inspiration: *locally orderless images* Koenderink & Van Doorn (1999)



- Previous work: "subdivide-and-disorder" strategy




Szummer & Picard (1997) SIFT: Lowe (1999, 2004) Gist: Torralba et al. (2003)

2

Spatial pyramid representation

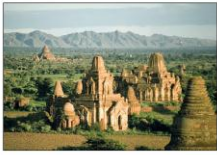
- Extension of a bag of features
- Locally orderless representation at several levels of resolution
- Based on *pyramid match kernels* Grauman & Darrell (2005)
 - Grauman & Darrell: build pyramid in feature space, discard spatial information
 - Our approach: build pyramid in image space, quantize feature space




level 0 level 1 level 2

3

Feature extraction




Weak features



Edge points at 2 scales and 8 orientations (vocabulary size 16)

Strong features



SIFT descriptors of 16x16 patches sampled on a regular grid, quantized to form visual vocabulary (size 200, 400)

5

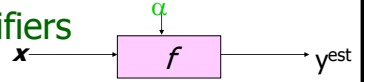
Support Vector Machines

Modified from the slides by Dr. Andrew W. Moore
<http://www.cs.cmu.edu/~awm/tutorials>

Copyright © 2001, 2003, Andrew W. Moore

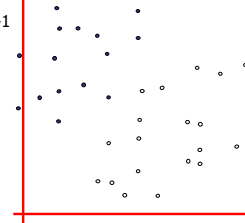
Nov 23rd, 2001

Linear Classifiers



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1

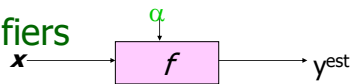


How would you classify this data?

Copyright © 2001, 2003, Andrew W. Moore

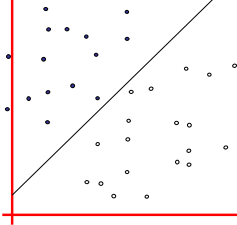
Support Vector Machines: Slide 14

Linear Classifiers



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1

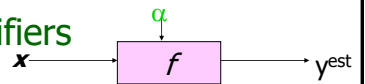


How would you classify this data?

Copyright © 2001, 2003, Andrew W. Moore

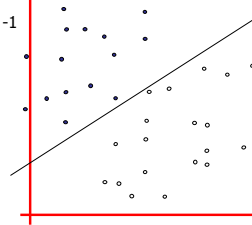
Support Vector Machines: Slide 15

Linear Classifiers



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1



How would you classify this data?

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 16

Linear Classifiers

α

\mathbf{x} → f → y_{est}

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

How would you classify this data?

Copyright © 2001, 2003, Andrew W. Moore Support Vector Machines: Slide 17

Linear Classifiers

α

\mathbf{x} → f → y_{est}

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

Any of these would be fine..

..but which is best?

Copyright © 2001, 2003, Andrew W. Moore Support Vector Machines: Slide 18

Classifier Margin

α

\mathbf{x} → f → y_{est}

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

Copyright © 2001, 2003, Andrew W. Moore Support Vector Machines: Slide 19

Maximum Margin

α

\mathbf{x} → f → y_{est}

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin.

This is the simplest kind of SVM (Called an LSVM)

Linear SVM

Copyright © 2001, 2003, Andrew W. Moore Support Vector Machines: Slide 20

Maximum Margin

α

$x \rightarrow f \rightarrow y_{est}$

$f(x, w, b) = \text{sign}(w \cdot x - b)$

- denotes +1
- denotes -1

Support Vectors are those datapoints that the margin pushes up against

The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin.

This is the simplest kind of SVM (Called an LSVM)

Linear SVM

Copyright © 2001, 2003, Andrew W. Moore Support Vector Machines: Slide 21

Why Maximum Margin?

- denotes +1
- denotes -1

Support Vectors are those datapoints that the margin pushes up against

1. Intuitively this feels safest.
2. If we've made a small error in the location of the boundary (it's been jolted in its perpendicular direction) this gives us least chance of causing a misclassification.
3. LOOCV is easy since the model is immune to removal of any non-support-vector datapoints.
4. There's some theory (using VC dimension) that is related to (but not the same as) the proposition that this is a good thing.
5. Empirically it works very very well.

Copyright © 2001, 2003, Andrew W. Moore Support Vector Machines: Slide 22

Nonlinear Kernel (I)

Example: SVM with Polynomial of Degree 2

Kernel: $K(\vec{x}_i, \vec{x}_j) = [\vec{x}_i \cdot \vec{x}_j + 1]^2$

plot by Bell SVM applet

Copyright © 2001, 2003, Andrew W. Moore Support Vector Machines: Slide 23

Nonlinear Kernel (II)

Example: SVM with RBF-Kernel

Kernel: $K(\vec{x}_i, \vec{x}_j) = \exp(-|\vec{x}_i - \vec{x}_j|^2 / \sigma^2)$

plot by Bell SVM applet

Copyright © 2001, 2003, Andrew W. Moore Support Vector Machines: Slide 24

Scene category dataset

Fei-Fei & Perona (2005), Oliva & Torralba (2001)
http://www-cvr.ai.uiuc.edu/ponce_grp/data

Multi-class classification results (100 training images per class)

Level	Weak features (vocabulary size: 16)		Strong features (vocabulary size: 200)	
	Single-level	Pyramid	Single-level	Pyramid
0 (1 × 1)	45.3 ± 0.5		72.2 ± 0.6	
1 (2 × 2)	53.6 ± 0.3	56.2 ± 0.6	77.9 ± 0.6	79.0 ± 0.5
2 (4 × 4)	61.7 ± 0.6	64.7 ± 0.7	79.4 ± 0.3	81.1 ± 0.3
3 (8 × 8)	63.3 ± 0.8	66.8 ± 0.6	77.2 ± 0.4	80.7 ± 0.3

Fei-Fei & Perona: 65.2%

Caltech101 dataset

Fei-Fei et al. (2004)
http://www.vision.caltech.edu/Image_Datasets/Caltech101/Caltech101.html

Multi-class classification results (30 training images per class)

Level	Weak features (16)		Strong features (200)	
	Single-level	Pyramid	Single-level	Pyramid
0	15.5 ± 0.9		41.2 ± 1.2	
1	31.4 ± 1.2	32.8 ± 1.3	55.9 ± 0.9	57.0 ± 0.8
2	47.2 ± 1.1	49.3 ± 1.4	63.6 ± 0.9	64.6 ± 0.8
3	52.2 ± 0.8	54.0 ± 1.1	60.3 ± 0.9	64.6 ± 0.7

Caltech 101

Caltech256

[Description] [Download] [Discussion] [Other Datasets]

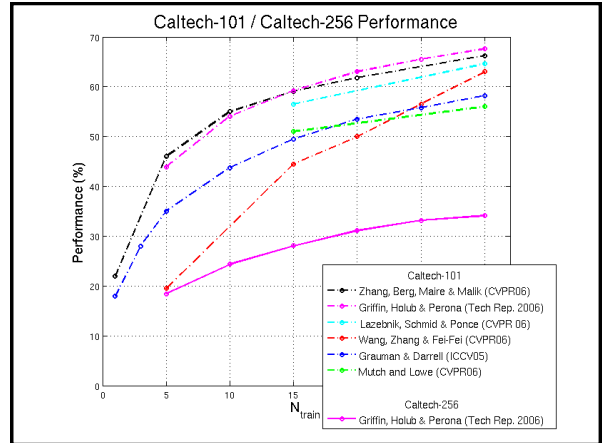
Description

Pictures of objects belonging to 101 categories. About 40 to 800 images per category. Most categories have about 50 images. Collected in September 2003 by Fei-Fei Li, Marco Andreetto, and Marc Aurelio Ranzato. The size of each image is roughlyly 300 pixels. We have carefully clicked outlines of each object in these pictures, these are included under the 'Annotations.tar'. There is also matlab script to view the annotations, 'show_annotations.m'.

How to use the dataset

Caltech-101: Drawbacks

- Smallest category size is 31 images: $N_{train} \leq 30$
- Too easy?
 - left-right aligned
 - Rotation artifacts
 - Soon will saturate performance



- Jump to Nicolas Pinto's slides. (page 29)



Papers

- A. Torralba. Contextual priming for object detection. IJCV 2003.
- A. Rabinovich, A. Vedaldi, C. Galleguillos, E. Wiewiora and S. Belongie. Objects in Context. ICCV 2007

33

Object Detection

Probabilistic Framework

$$P(O | \mathbf{v}) = \frac{P(\mathbf{v} | O)}{P(\mathbf{v})} P(O) \quad (\text{Single Object Likelihood})$$

Object presence at a particular location/scale
Given all image features (local/object and scene/context)

$$\mathbf{v} = \mathbf{v}_{\text{Local}} + \mathbf{v}_{\text{Contextual}}$$

34

Contextual Reasoning

2D Reasoning

2,5D / 3D Reasoning

Scene Centered

Object Centered

Surface orientations w.r.t. camera



Contextual priming for object detection



Objects in Context



Geometric context from a single image.

35

Preview: Contextual Priming for Object Detection



Input test image

36

Preview: Contextual Priming for Object Detection

Correlate with many filters

37

Preview: Contextual Priming for Object Detection

Using previously collected statistics about filter output *predict* information about objects

38

Preview: Contextual Priming for Object Detection

Predict information about objects

Where I can find the objects easily?

people chair car
Which objects do I expect to see?

How large objects do I expect to see?

39

Contextual Priming for Object Detection: Probabilistic Framework

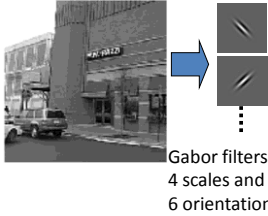
$$P(O | \mathbf{v}) = \frac{P(O, \mathbf{v})}{P(\mathbf{v})} = \frac{P(\mathbf{v}_L | O, \mathbf{v}_C)}{P(\mathbf{v}_L | \mathbf{v}_C)} P(O | \mathbf{v}_C)$$

Local measurements (a lot in the literature)

Contextual features

40

Contextual Priming for Object Detection: Contextual Features

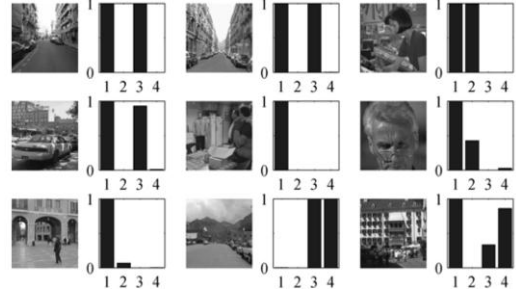


Use **PCA** on filter output images to reduce the number of features (< 64)

Use **Mixture of Gaussians** to model the probabilities. (Other alternatives include KNN, parzen window, logistic regression, etc)

41

Contextual Priming for Object Detection: Object Priming Results



(o_1 =people, o_2 =furniture, o_3 =vehicles and o_4 =trees)

42

Contextual Priming for Object Detection: FOCUS of Attention Results



Heads

43

Contextual Priming for Object Detection: Conclusions

- Proves the relation btw low level features and scene/context
- Can be seen as a computational evidence for the (possible) existence of low-level feature based biological attention mechanisms
- Also a warning: Whether an object recognition system understands the object or works by lots of features.

44

Preview: Objects in Context



Input test image

45

Preview: Objects in Context



Do segmentation on the image

46

Preview: Objects in Context



Do classification (find label probabilities) in each segment only with local info

47

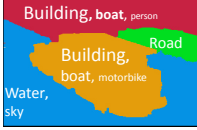
Preview: Objects in Context



Most consistent labeling according to *object co-occurrences* & local label probabilities.

48

Objects in Context: Local Categorization




- Extract random patches on zero-padded segments
- Calculate SIFT descriptors
- Use BoF:
 - Training:
 - Cluster patches in training (Hier. K-means, $K=10x3$)
 - Histogram of words in each segment
 - NN classifier (returns a sorted list of categories)

Each segment is classified independently

49

Objects in Context: Contextual Refinement



Contextual model based on co-occurrences
Try to find the most consistent labeling with **high posterior probability** and **high mean pairwise interaction**.
Use CRF for this purpose.

$$p(c_1 \dots c_k | S_1 \dots S_k) = \frac{B(c_1 \dots c_k) \prod_{i=1}^k A(i)}{Z(\phi, S_1 \dots S_k)}$$

$B(c_1 \dots c_k) = \exp\left(\sum_{i,j=1}^k \phi(c_i, c_j)\right)$
 Mean interaction of all label pairs
 $\Phi(i,j)$ is basically the observed label co-occurrences in training set.

$A(i) = p(c_i | S_i)$
 Independent
 segment classification

50

Objects in Context: Learning Context

Using labeled image datasets (MSRC, PASCAL)

Using labeled text based data (Google Sets): Contains list of related items
→ A large set turns out to be useless! (anything is related)

51

Objects in Context: Results

	No Context	Google Sets	Using Training
MSRC	45.0%	58.1%	68.4%
PASCAL	61.8%	63.4%	74.2%

Table 1. Average Categorization Accuracy.

52

“Objects in Context” – Limitations: Context modeling

Segmentation

Categorization without context
Local information only

With **co-occurrence** context
Means: $P(\text{person,dog}) > P(\text{person,cow})$

(Bonus Q: How did it handle the background?)

53

“Objects in Context” – Limitations: Context modeling

Segmentation

Categorization without context
Local information only

With **co-occurrence** context
 $P(\text{person,horse}) > P(\text{person,dog})$

But why? Isn't it only a dataset bias? We have seen in the previous example that $P(\text{person,dog})$ is common too.

54

“Objects in Context” Object-Object or Stuff-Object ?

MSRC training data

building	75	18	29	33	4	9	7	18	10	2	1	43	1	9	4
grass	95	36	23	15	39	14	7	7	3	1	4	15	2	5	6
tree	38	65	6	43	5	12	4	4	1	2	1	19	11	8	
cow	23	97	23	7	4										
sheep	15		15		1										
sky	39	43	4	06	19	18	4	3	5	4	2	5	8	11	
aeroplane	14	47		19	18										
water	7	10	4	1	15	43	4	1	7	5	3	8	12		
face	19	4	3	1	1	20						19	1		
bike	59	9				1	15					19	1		
flower	2	1	1	1	4	7						14	3	1	
sign	2	2	5												
bird	1	1	1	1	1	1									
book	4	1	1	1	1	1									
chair	4	1	1	1	1	1									
road	43	15	19	2	25	5	8	7	19	12	1	3	3	7	10
cat	1	2										19	13	1	
dog	5	8	11	5	6	20	1	1	1	1	1	5	8	32	2
body	1	1	1	1	1	1									
boat															
building															
grass															
sheep															
cow															
aeroplane															
face															
bike															
flower															
bird															
book															
chair															
cat															
dog															
boat															

Stuff-like

Stuff

Labels with high co-occurrences with other labels

Looks like “background” stuff – object (such as water-boat) does help rather than “foreground” object co-occurrences (such as person-horse) [but still car-person-motorbike is useful in PASCAL]

55

“Objects in Context” – Limitations: Segmentation

- Too good: A few or many? How to select a good segmentation in multiple segmentations?
- Can make object recognition & contextual reasoning (due to stuff detection) much easier.


56

“Objects in Context” - Limitations

- No cue by **unknown objects**
- No **spatial relationship** reasoning
- Object detection part heavily depends on **good segmentations**
- Improvements using object co-occurrences are demonstrated with images where **many labels are already correct**. → How good is the model?

57

Contextual Priming vs. Objects in Context

<p>Scene->Object</p> <p>Simpler training data (only target object's labels are enough)</p> <p>Scene information is view-dependent (due to gist)</p>  <p>Object detector independent</p>	<p>{Object,Stuff} <-> {Object,Stuff}</p> <p>May need huge amount of labeled data</p> <p>Can be more generic than scene->object with a very good model</p> <p>Contextual model is object detector independent, in theory. But:</p> <p>+ use segmentation → easier to detect stuff</p> <p>- uses segmentation → can be unreliable</p>
---	---

58

Microsoft Research

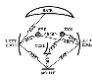
Finding the weakest link in person detectors

Devi Parikh
TTI, Chicago

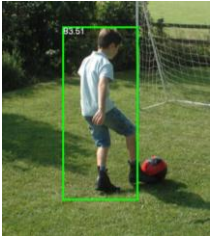
Larry Zitnick
Microsoft Research

Object recognition

We've come a long way...

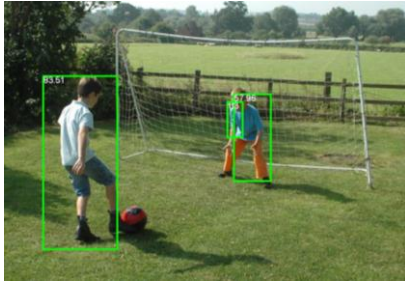


Fischler and Eilschlager, 1973

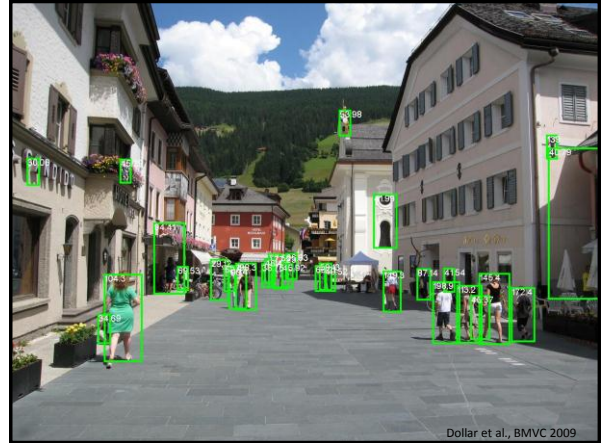


Dollar et al., BMVC 2009

Still a ways to go...



Dollar et al., BMVC 2009



Dollar et al., BMVC 2009

Still a ways to go...



Dollar et al., BMVC 2009

Part-based person detector

4 main components:

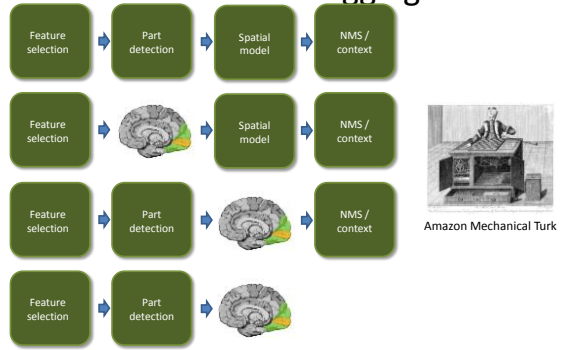


How can we help?

- Humans supply training data...
 - 100,000s labeled images
- We design the algorithms.
 - Going on 40 years.
- **Can we use humans to debug?**



Human debugging



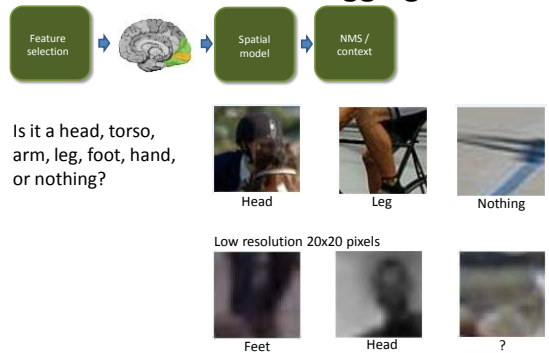
Human performance

- Humans ~90% average precision
- Machines ~46% average precision



PASCAL VOC dataset

Human debugging



Part detections

Humans Machine

P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE Trans. on PAMI*, 32:1627–1645, 2010.

Part detections

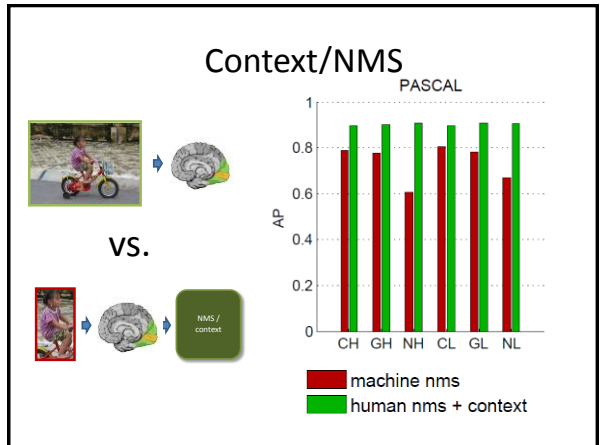
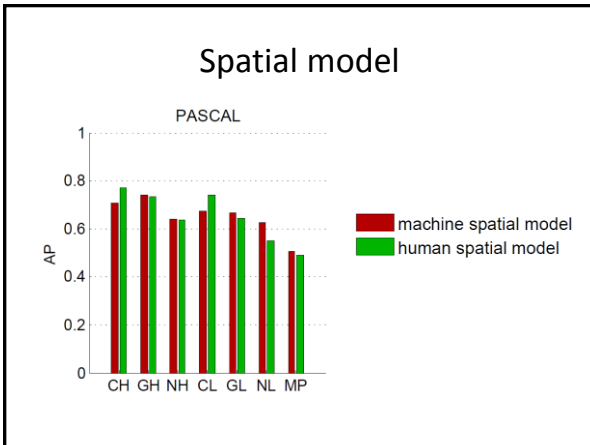
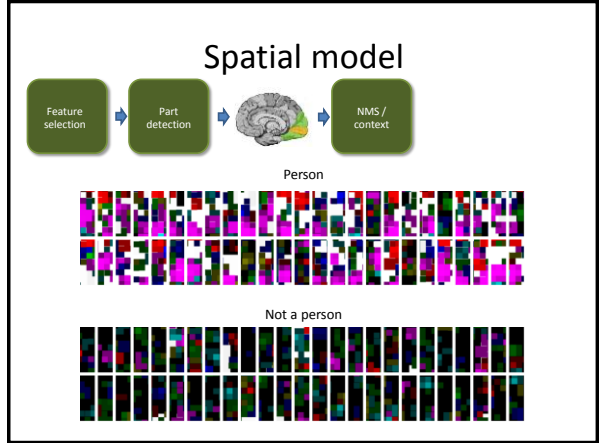
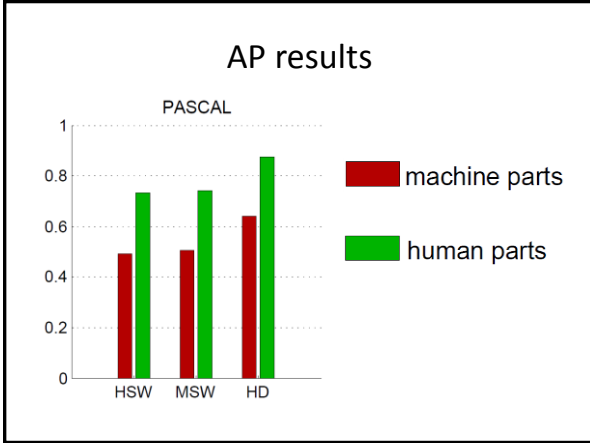
Humans Machine

Part detections

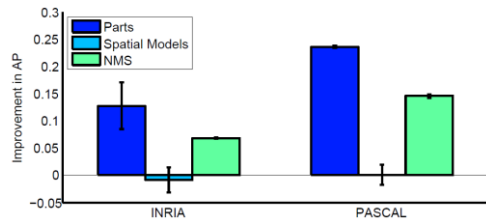
Humans Machine

Part detections

Humans Machine



Conclusion



http://www.ted.com/talks/lang/eng/pawan_sinha_on_how_brains_learn_to_see.html 7:00min

