

Machine Learning for Computer Vision

21 October, 2013
MVA – ENS Cachan



Rapid Object Detection with Deformable Part Models

Iasonas Kokkinos

iasonas.kokkinos@ecp.fr

Center for Visual Computing
Ecole Centrale Paris

Galen Group
INRIA-Saclay

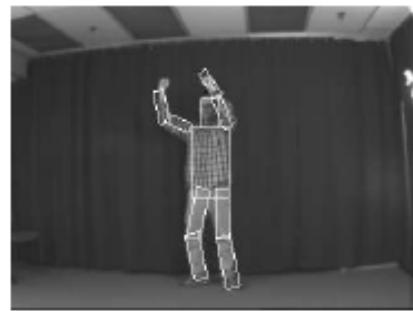
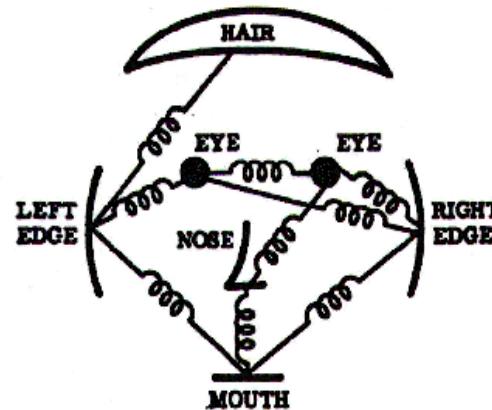
Perceptual Organization



Perceptual Organization



Object Detection with Deformable Part Models



P.F. Felzenszwalb, and D. Huttenlocher, *Pictorial Structures for Object Recognition*, IJCV 2005

P.F. Felzenszwalb, and D. Huttenlocher, *Distance Transforms of Sampled Functions*, 2004

P.F. Felzenszwalb, R. B. Girshick, and D. A. McAllester, *Cascade object detection with DPMs* CVPR 2010

Deformable Part Models (DPMs)

$$\mathbf{S}(\mathbf{x}) = \sum_{p=1}^P U_p(x_p) + B_p(x, x_p)$$

Local appearance

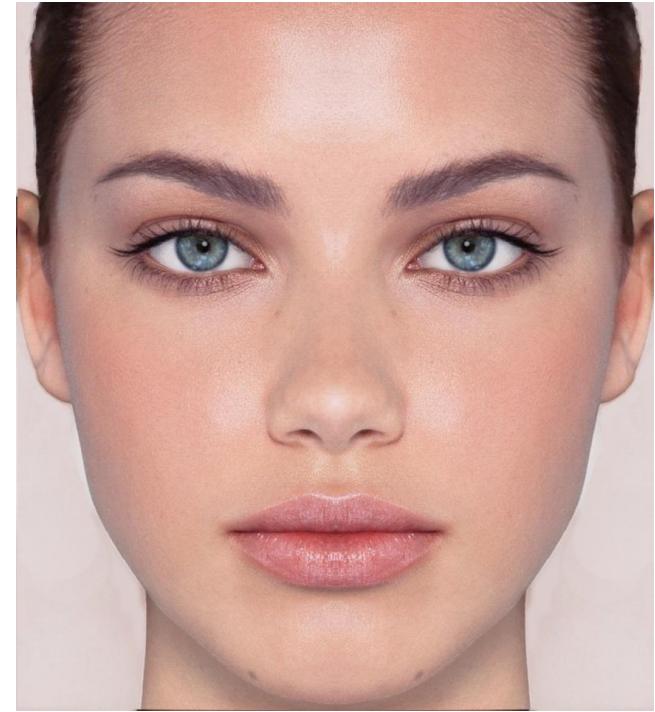
$$U_p(x_p) = \langle w_p, H(x_p) \rangle$$

Pairwise compatibility

$$B_p(x, x_p) = -(h - h_p - \hat{h}_p)^2 \eta - (v - v_p - \hat{v}_p)^2 \nu$$

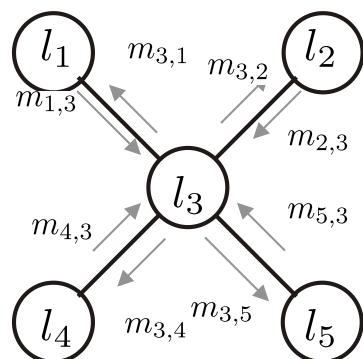
Object score

$$S(x) = \sum_{p=1}^P \max_{x'} [U_p(x') + B_p(x, x')]$$



$$\mu_{i \rightarrow j}(X_j) = \max_{X_i} \Phi_i(X_i) \Psi_{i,j}(X_i, X_j) \prod_{k \in \mathcal{N}(i) \setminus j} \mu_{k \rightarrow j}(X_i)$$

$$B_j(X_j) = \max P(X_j)$$



$$\Phi_1(l)$$



$$\Phi_2(l)$$



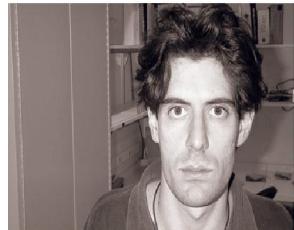
$$\Phi_3(l)$$



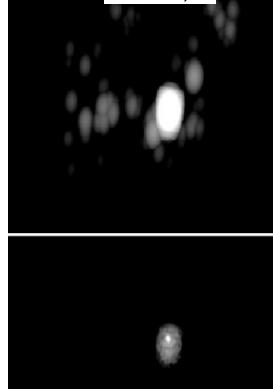
$$\Phi_4(l)$$



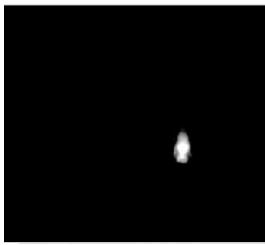
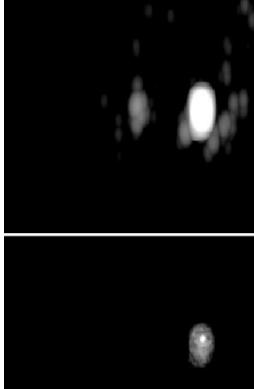
$$\Phi_5(l)$$



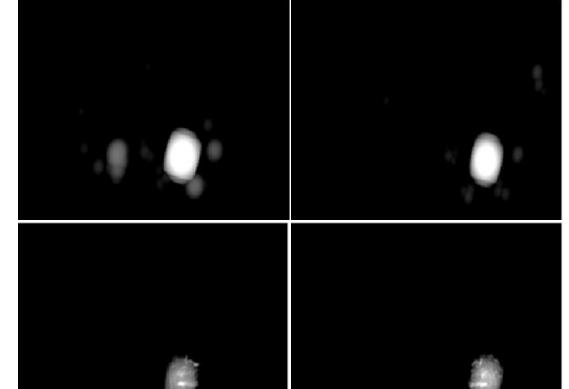
$$m_{1,3}$$



$$m_{2,3}$$



$$m_{3,4}$$



$$B_3(l)$$

$$B_1(l)$$

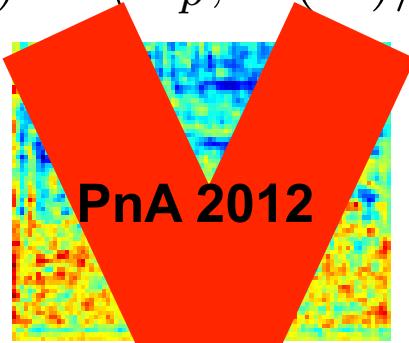
$$B_2(l)$$

$$B_4(l)$$

$$B_5(l)$$

Object detection with Deformable Part Models (DPMs)

$$U_p(x') = \langle \mathbf{w}_p, \mathbf{H}(x') \rangle$$



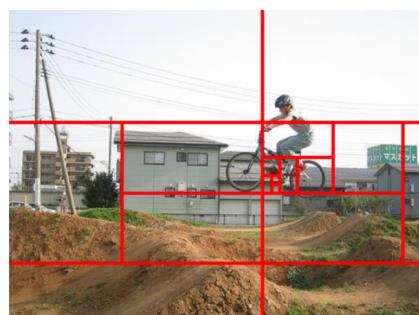
$$\max_{x'} [U_p(x') + B_p(x, x')]$$

$$p = 1$$

⋮

DTBB, NIPS 11

$$p = P$$



$$S(x) = \sum_{p=1}^P \max_{x'} [U_p(x') + B_p(x, x')]$$

Object Detection problem

$$S(x) = \sum_{p=1}^P \max_{x'} [U_p(x') + B_p(x, x')]$$

Threshold-based detection: recover all locations above threshold

$$M^\theta = \{x : S(x) \geq \theta\}$$

1-best detection: recover top-scoring candidate

$$M^* = \{\arg \max_x S(x)\}$$

Naïve implementation (Max-Product) : $O(P|x|^2)$

Generalized Distance Transforms (GDT): $O(P|x|)$

Our work (best-case): $O(|M|P \log_2 |x|)$

Key idea: avoid the exact evaluation of $S(x)$

Previous works on acceleration

Sparse image representations:

Y. Chen, L. Zhu, A. Yuille. Rapid inference on a novel and/or graph for object detection, NIPS 2007

I. Kokkinos and A. Yuille. HOP: Hierarchical Object Parsing, CVPR, 2009

Dense image representations:

B. Sapp, A. Toshev, and B. Taskar. Cascaded models for articulated pose estimation, ECCV, 2010

M. Pedersoli, A. Vedaldi, and J. Gonzalez. A coarse-to-fine approach for object detection, CVPR 2011

P F Felzenszwalb R B Girshick and D A McAllester Cascade object detection with DPMs CVPR 2010

‘Monolithic’ object models:

F. Fleuret and D. Geman. Coarse-to-fine face detection, IJCV 2001

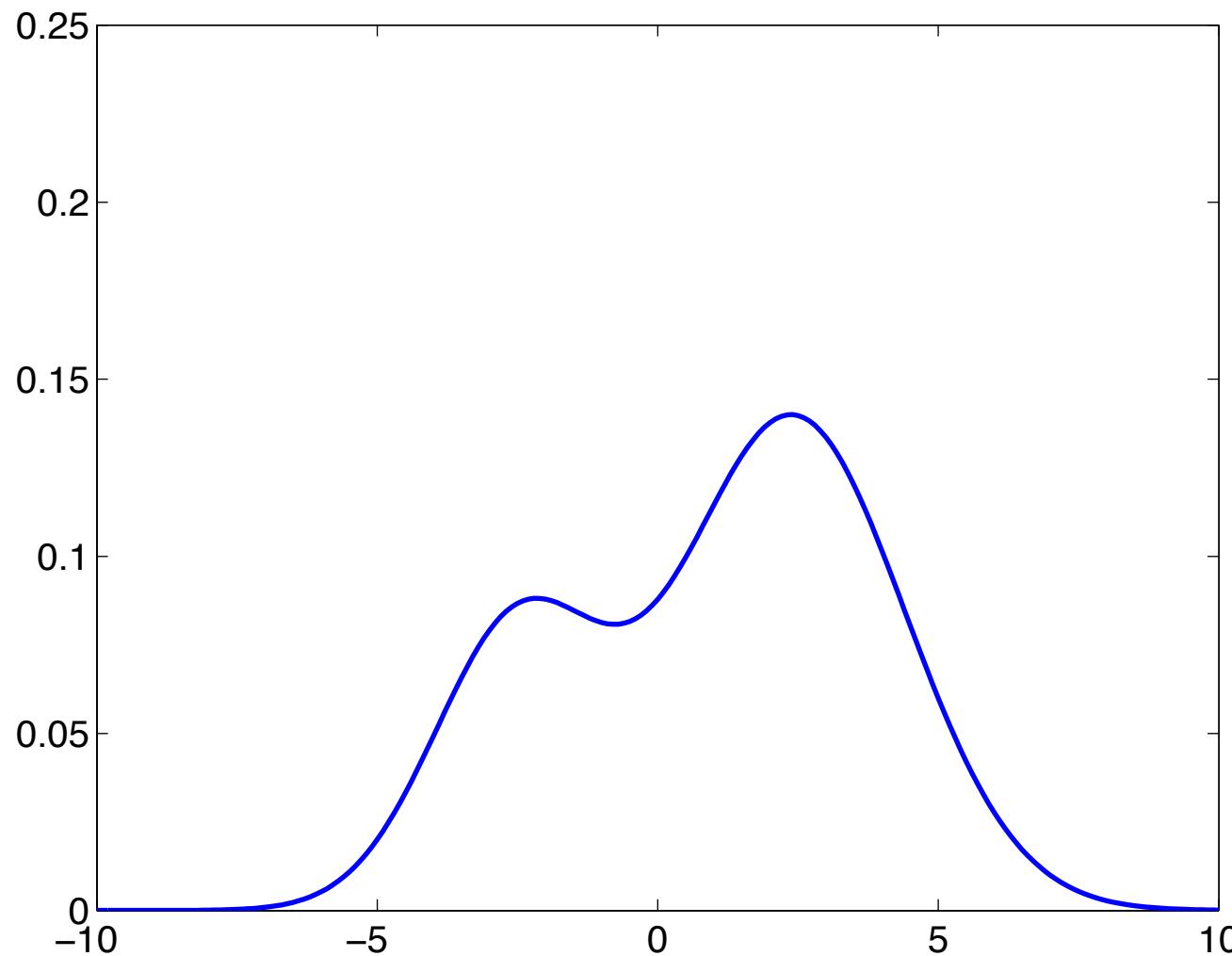
P. Viola and M. Jones. Robust Real-time Object Detection, IJCV 2004

C. Lampert , M. Blaschko and T. Hoffman. Beyond Sliding Windows- ESS, CVPR 2008

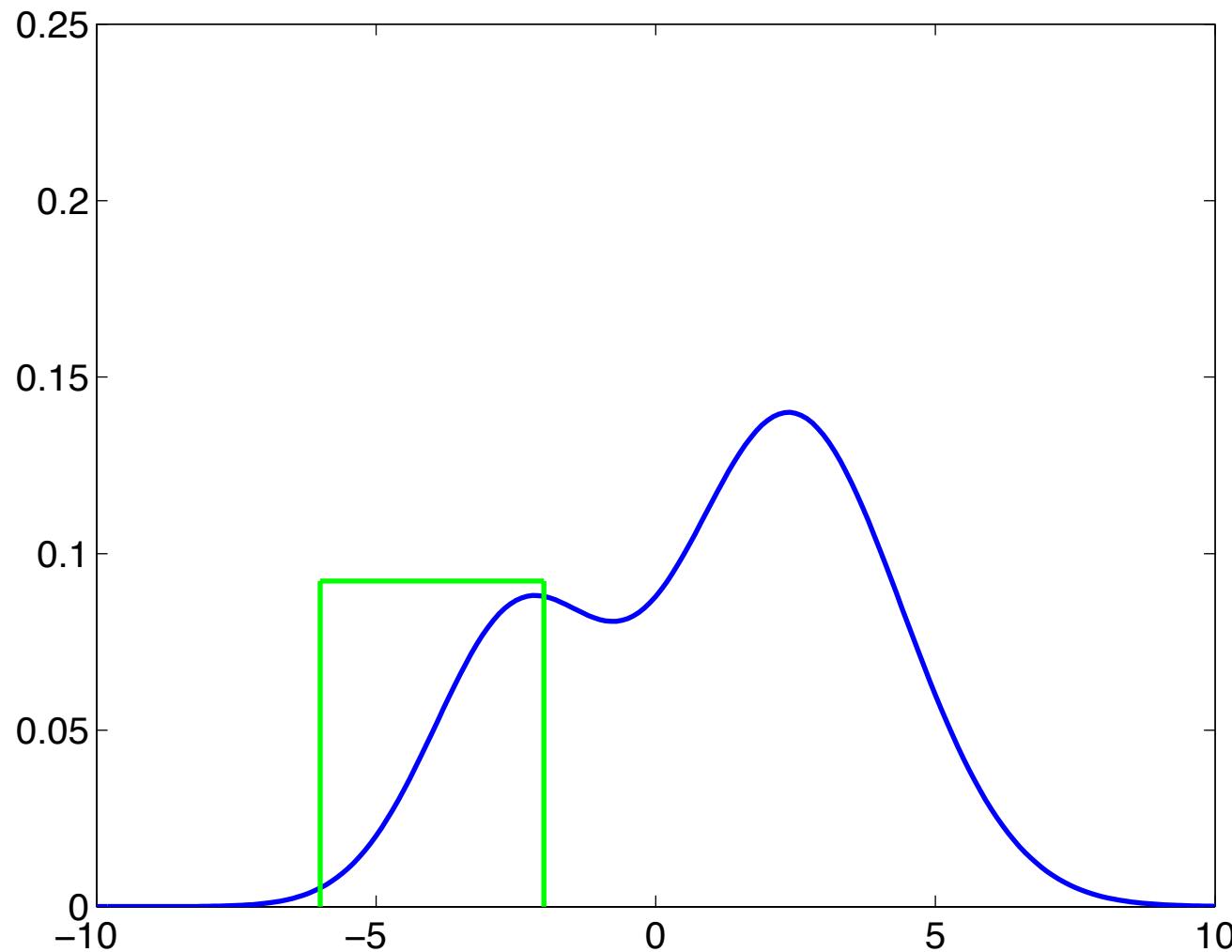
T. Breuel. Fast recognition using adaptive subdivisions of transformation space, CVPR 1992

Branch and bound

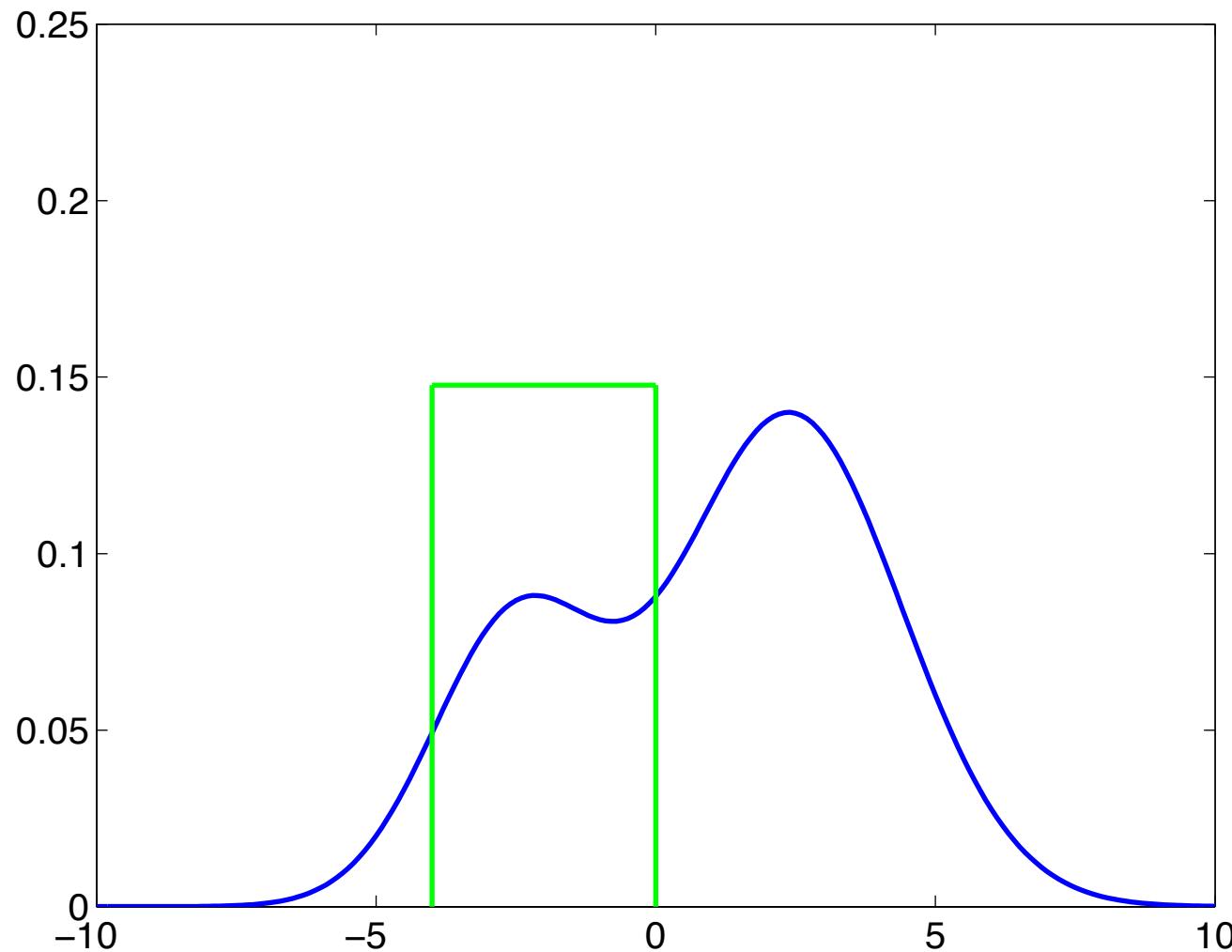
$$\max_{x \in X} f(x)$$



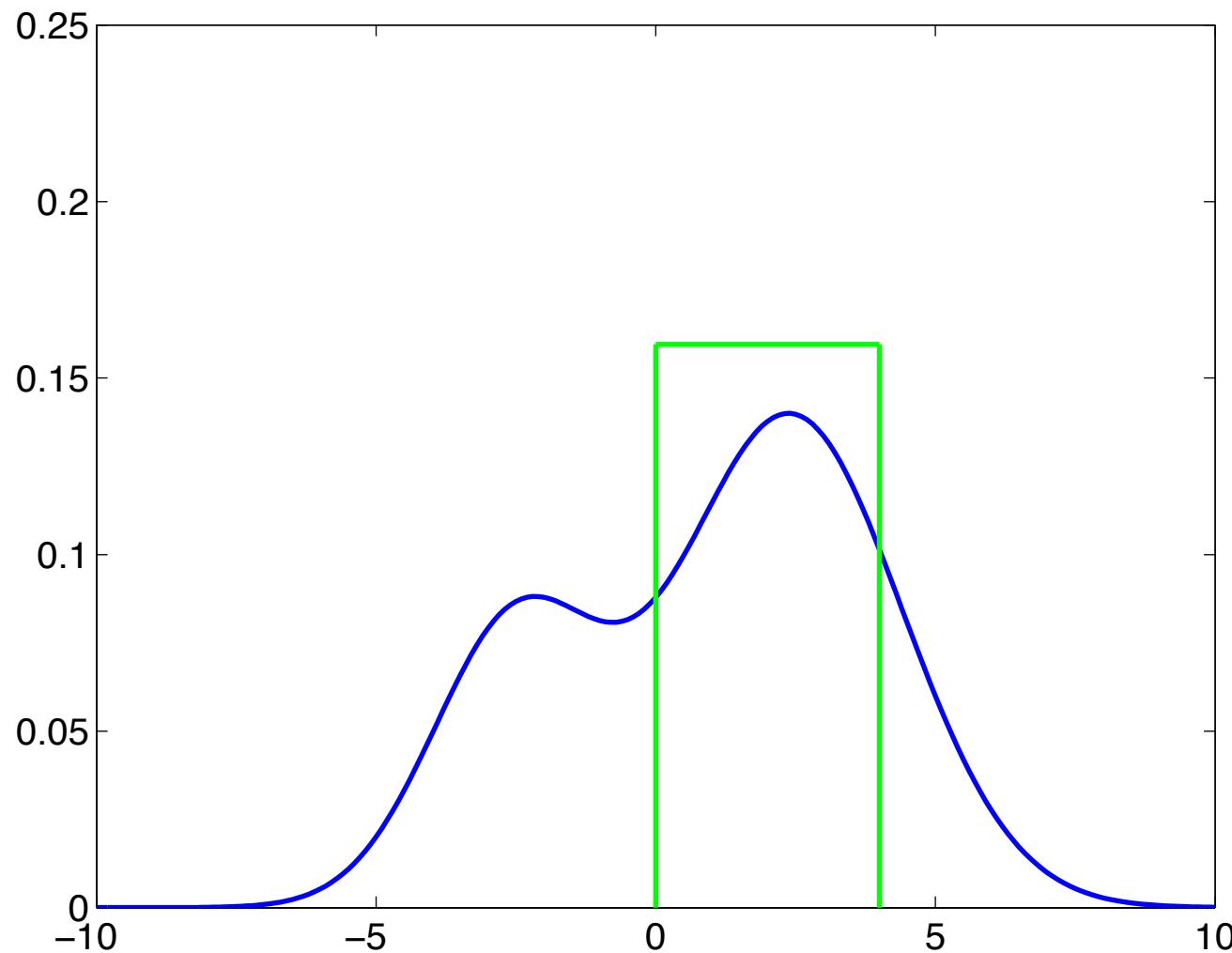
Bounding function $\bar{f}(X) \geq \max_{x \in X} f(x), \quad \bar{f}(\{x\}) = f(x)$



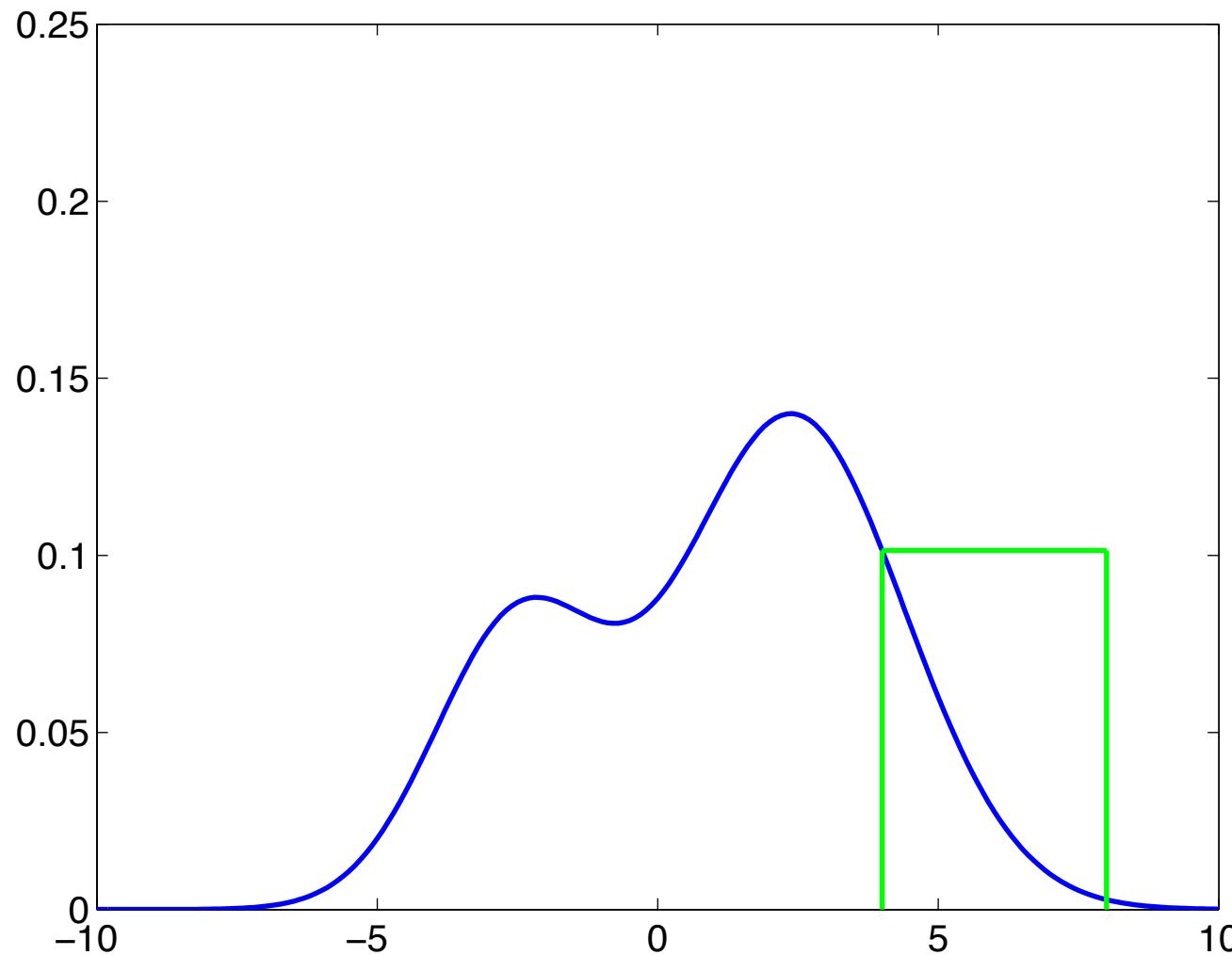
Bounding function $\bar{f}(X) \geq \max_{x \in X} f(x), \quad \bar{f}(\{x\}) = f(x)$



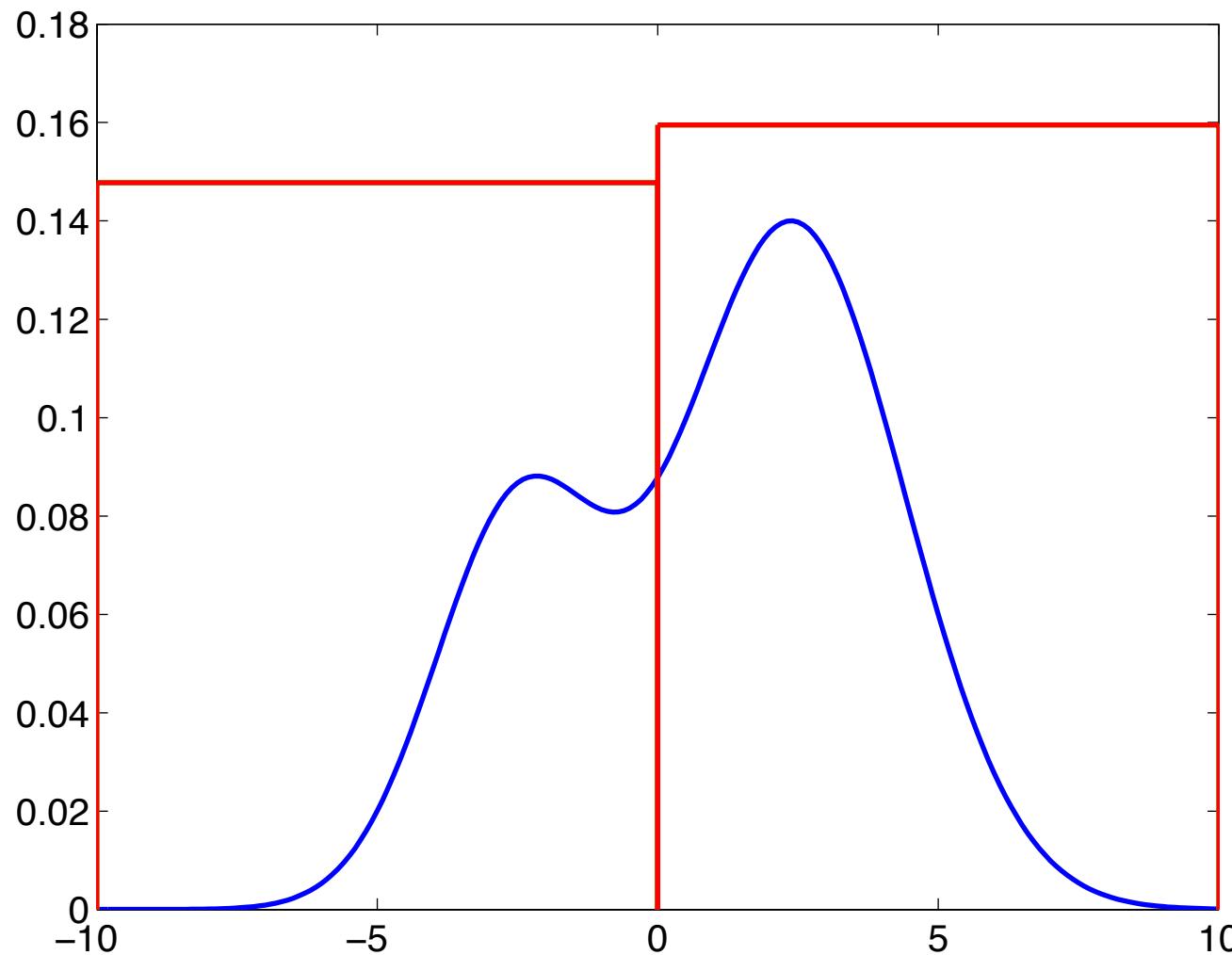
Bounding function $\bar{f}(X) \geq \max_{x \in X} f(x), \quad \bar{f}(\{x\}) = f(x)$



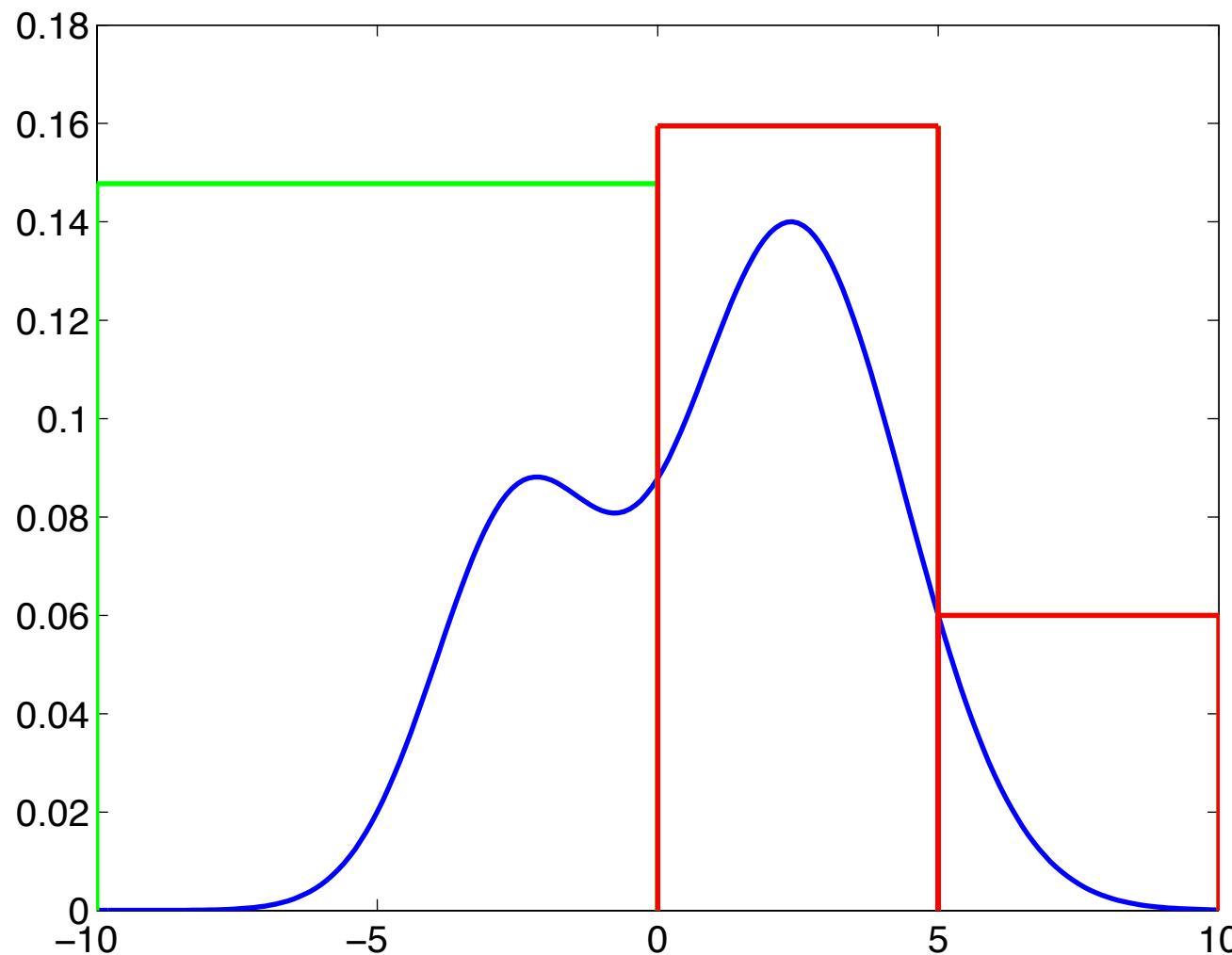
Bounding function $\bar{f}(X) \geq \max_{x \in X} f(x), \quad \bar{f}(\{x\}) = f(x)$



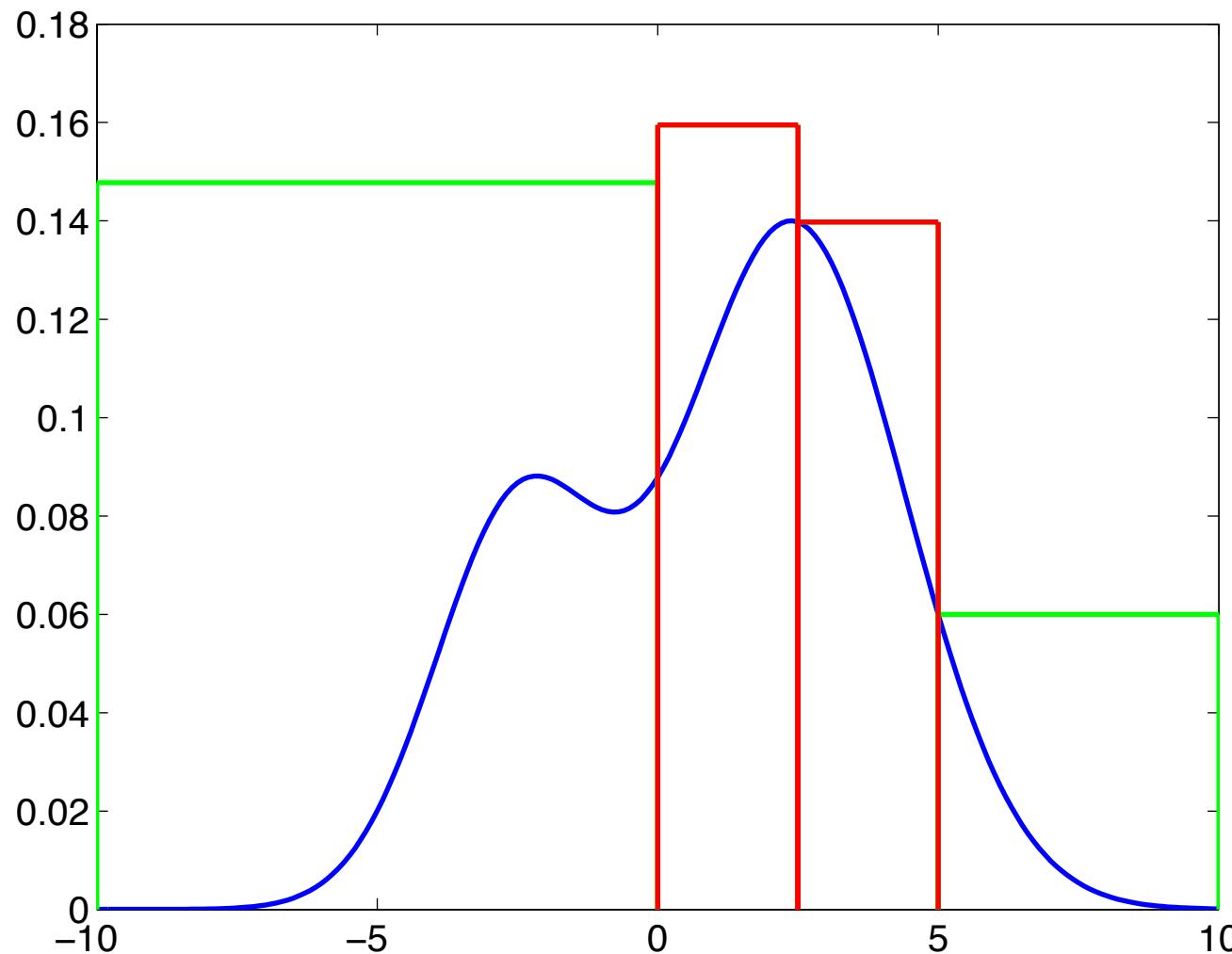
Branch-and-bound: all you need is bounds



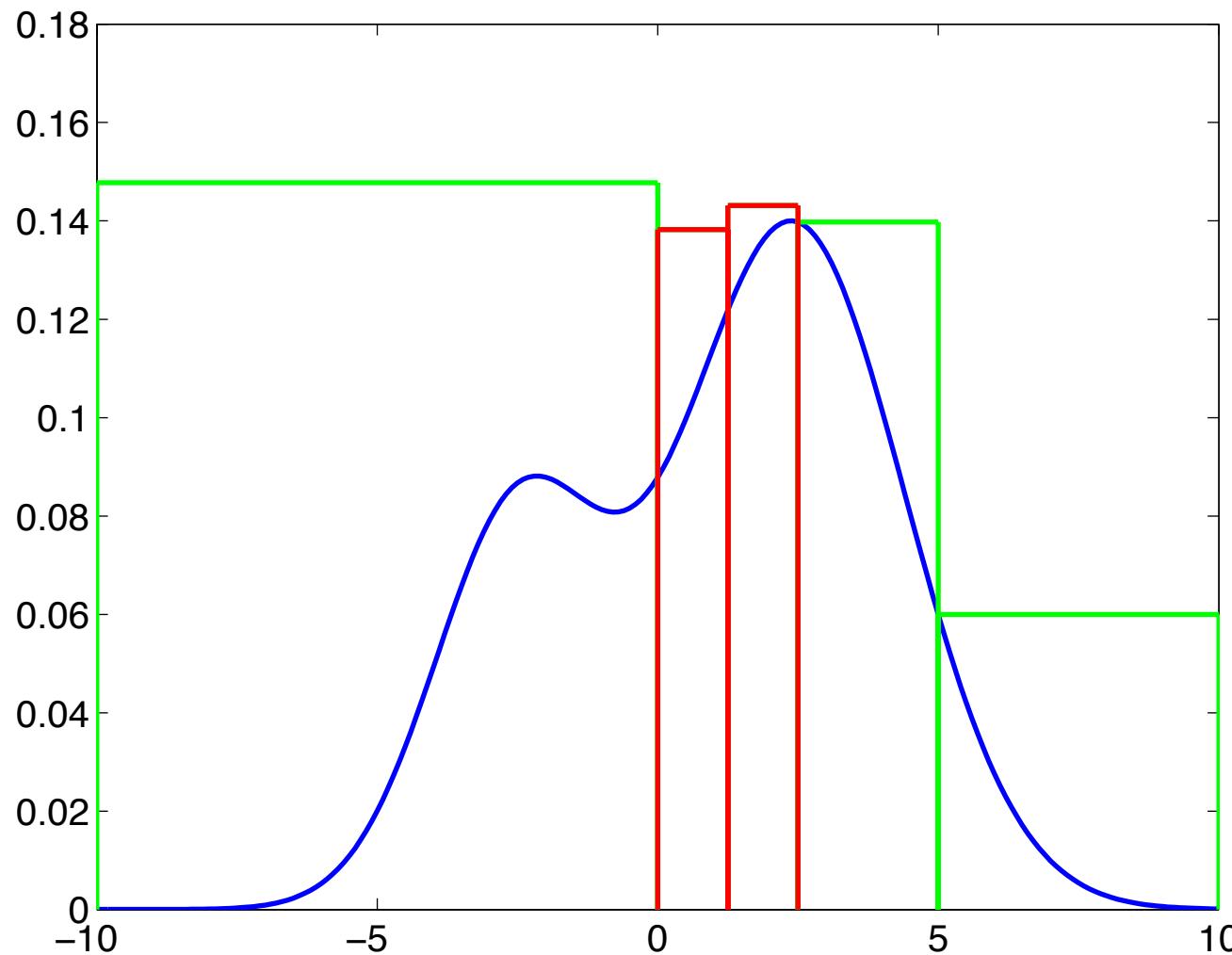
Branch-and-bound: all you need is bounds



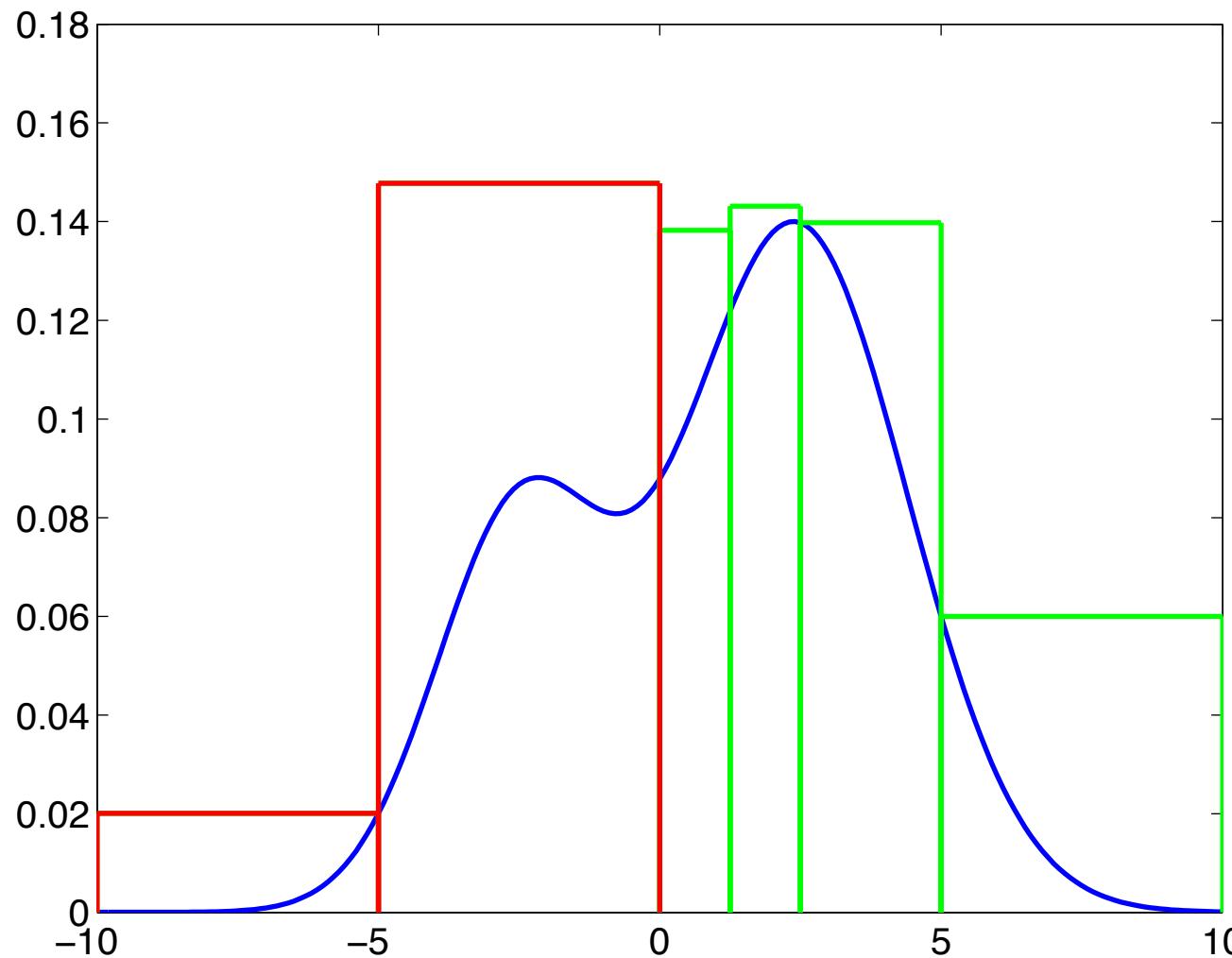
Branch-and-bound: all you need is bounds



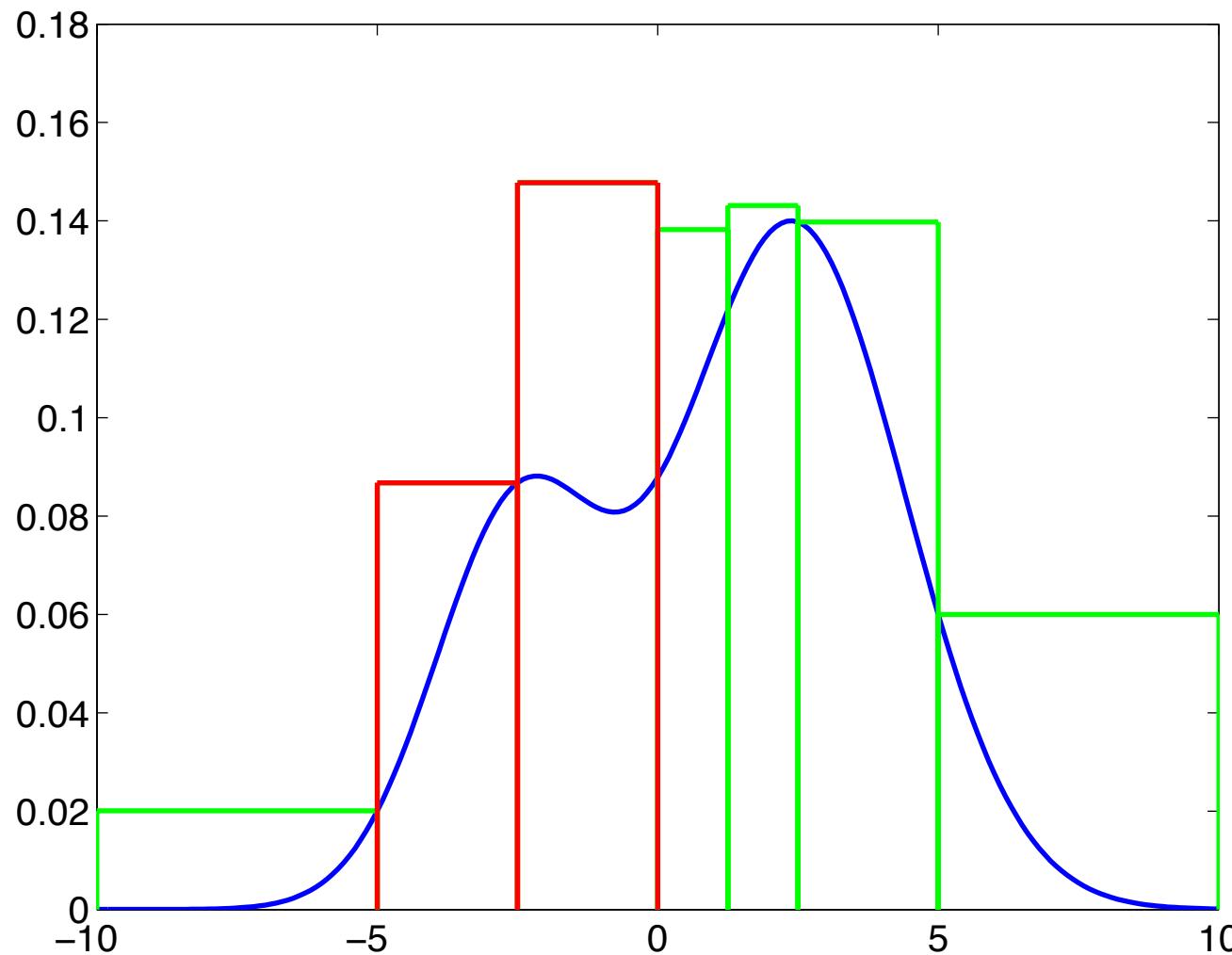
Branch-and-bound: all you need is bounds



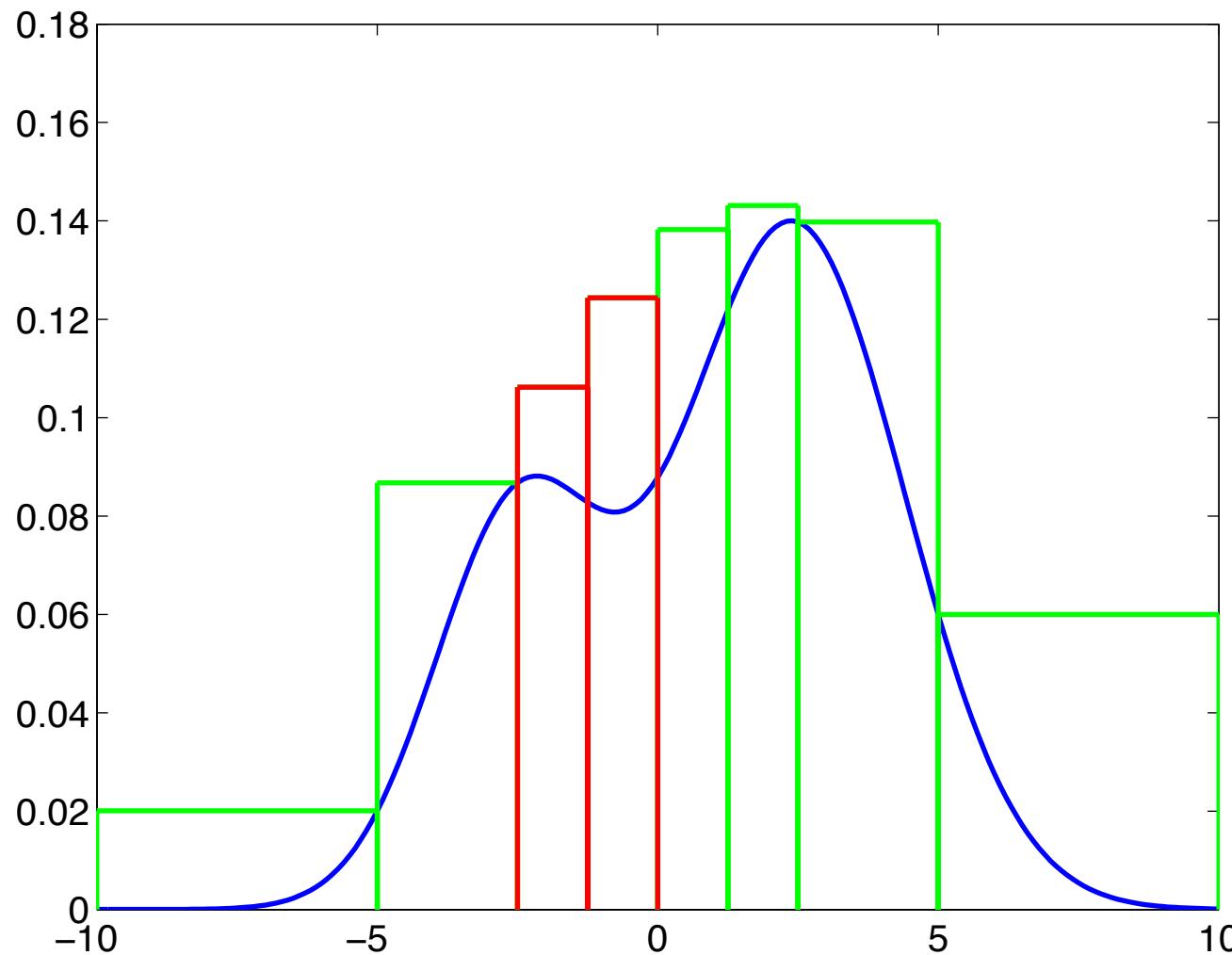
Branch-and-bound: all you need is bounds



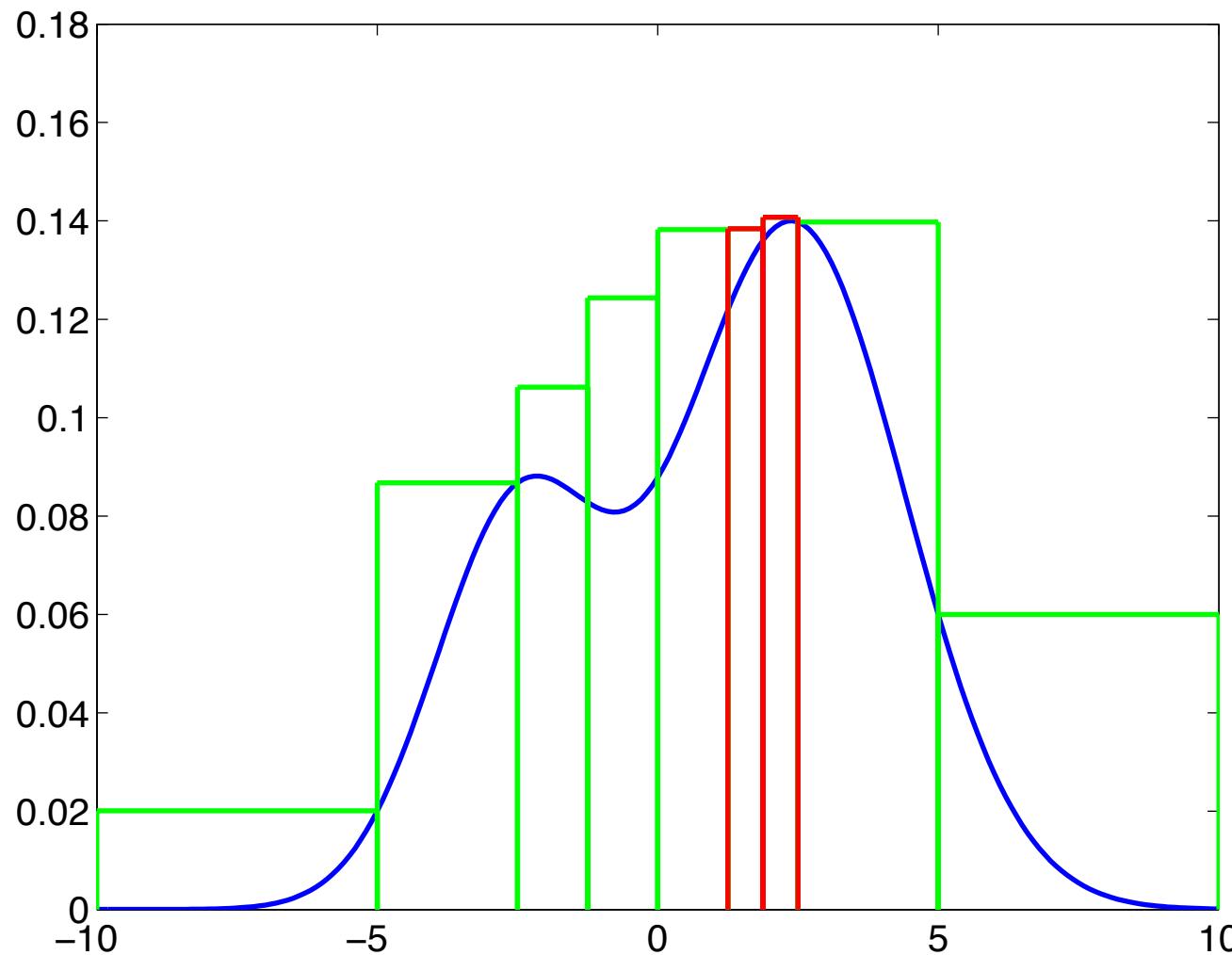
Branch-and-bound: all you need is bounds



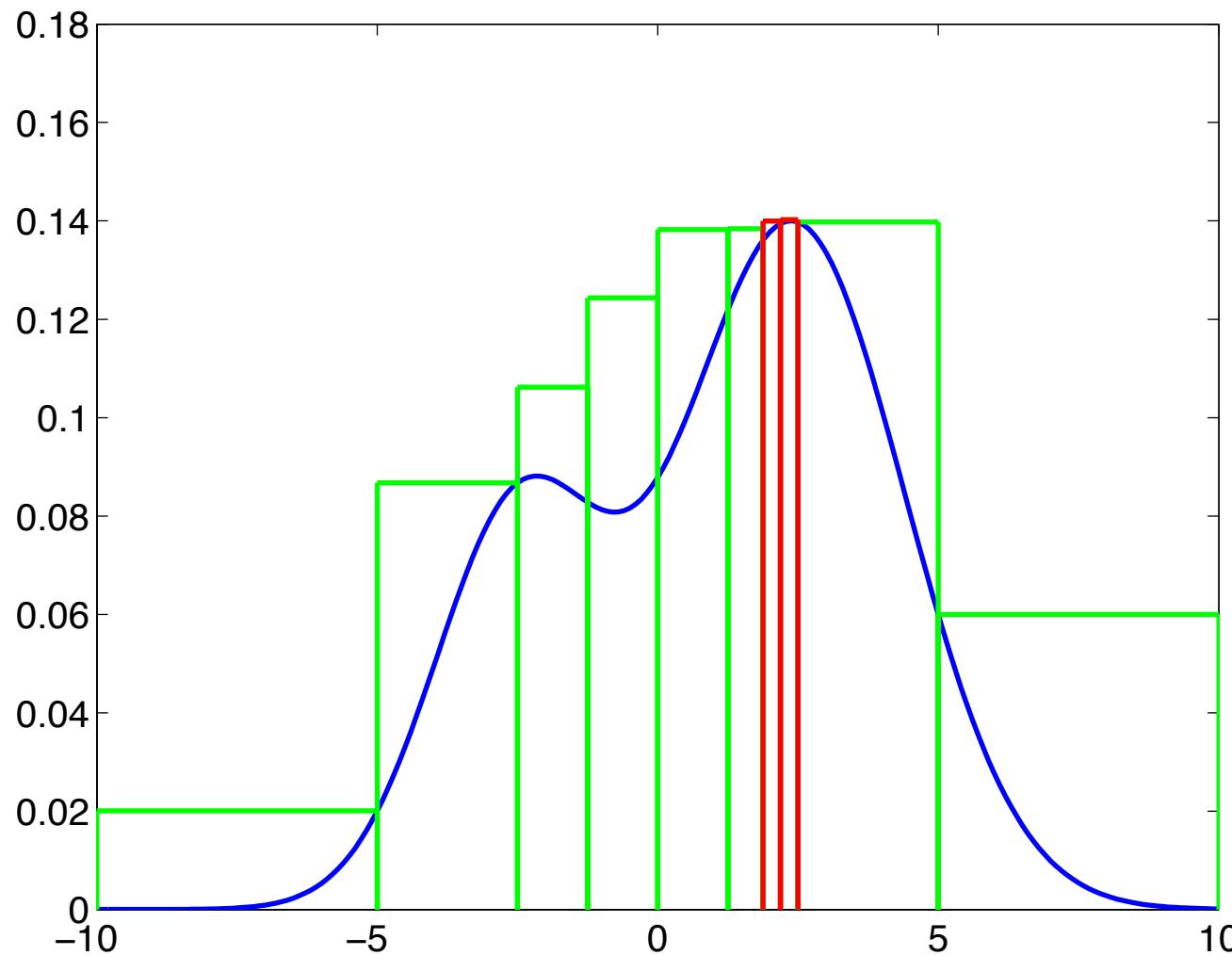
Branch-and-bound: all you need is bounds



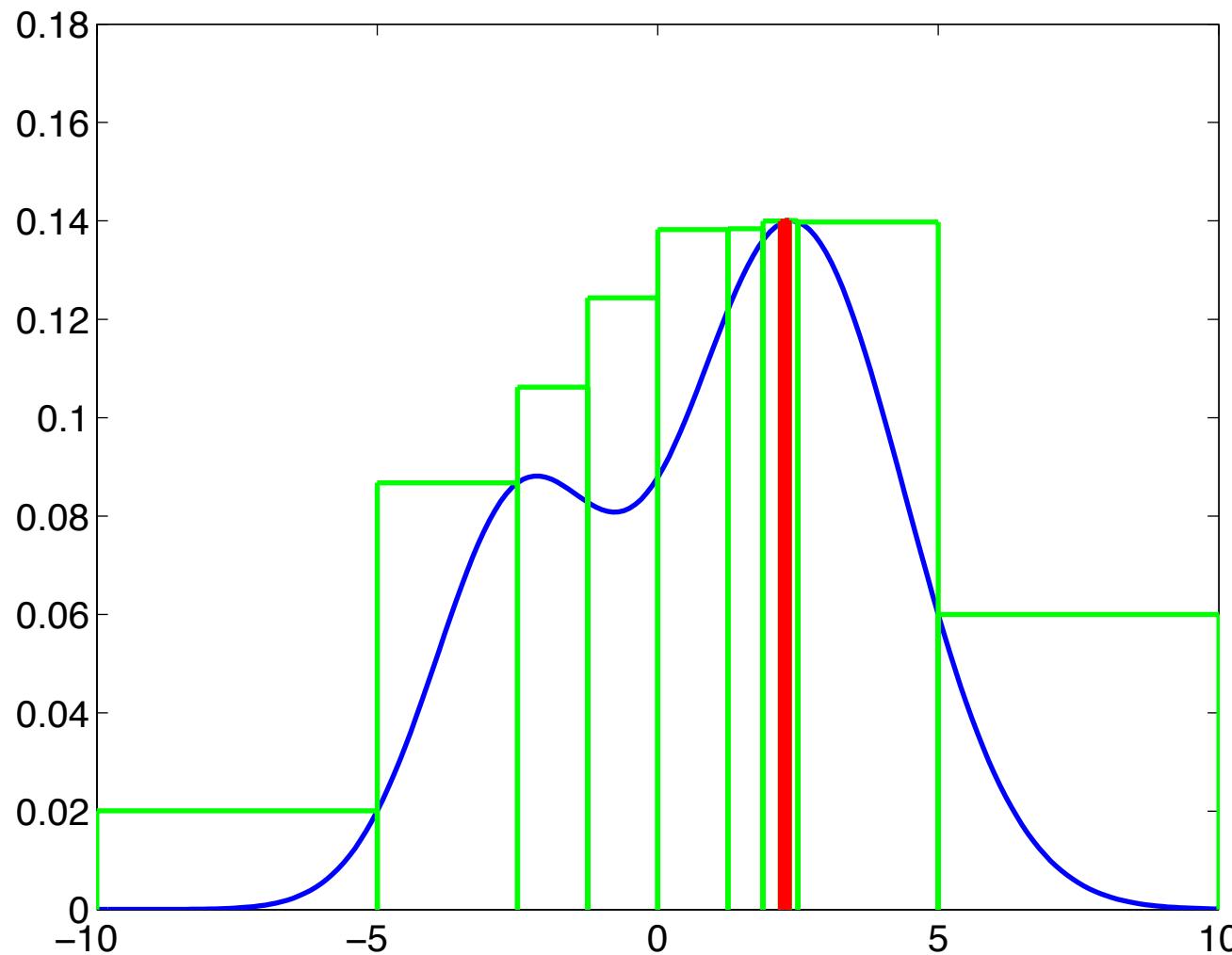
Branch-and-bound: all you need is bounds



Branch-and-bound: all you need is bounds



Branch-and-bound: all you need is bounds

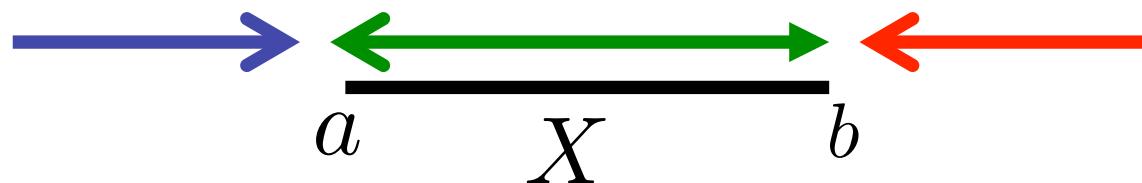


Bounding a mixture-of-gaussians

- Property: $\max_{x \in X} h(x) + g(x) \leq \max_{x \in X} h(x) + \max_{x \in X} g(x)$
- Function: $f(x) = \pi_1 N(x; \mu_1, \sigma_1) + \pi_2 N(x; \mu_2, \sigma_2)$

$$\begin{aligned} \max_{x \in X} f(x) &\leq \max_{x \in X} [\pi_1 N(x; \mu_1, \sigma_1)] + \max_{x \in X} [\pi_2 N(x; \mu_2, \sigma_2)] \\ &= \pi_1 N(d(X, \mu_1, \sigma_1); 0, 1) + \pi_2 N(d(X, \mu_2, \sigma_2); 0, 1) \end{aligned}$$

$$d(X, \mu, \sigma) = \begin{cases} 0 & a \leq \mu \leq b \\ \frac{1}{\sigma^2}(\mu - a)^2 & \mu \leq a \\ \frac{1}{\sigma^2}(\mu - b)^2 & b < \mu \end{cases}$$

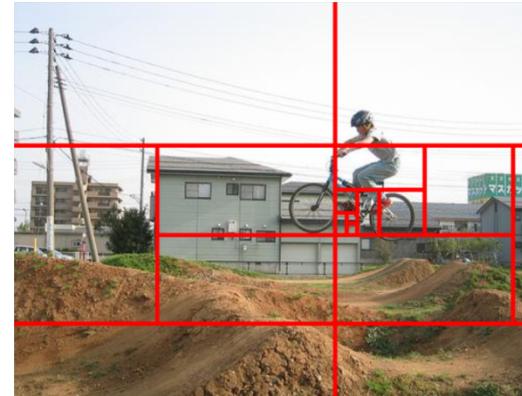
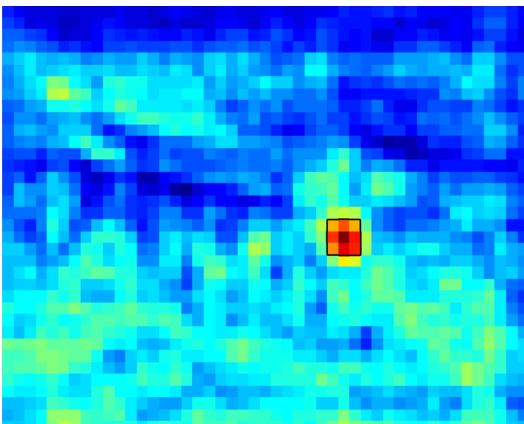


This work: BB for Deformable Part Models

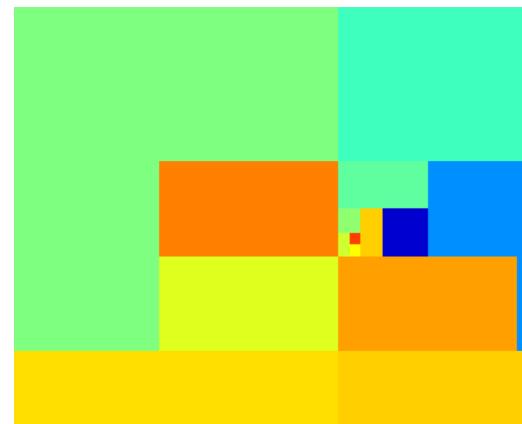
Input & Detection result



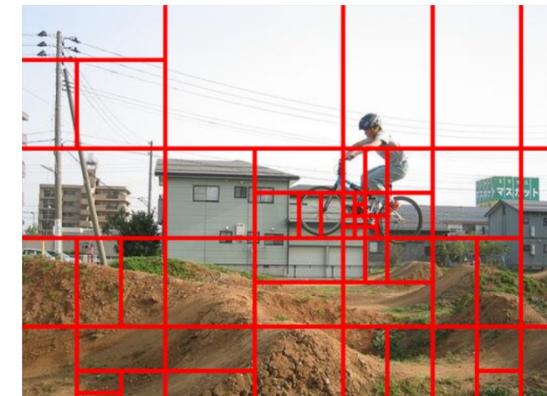
Detector score $S(x)$



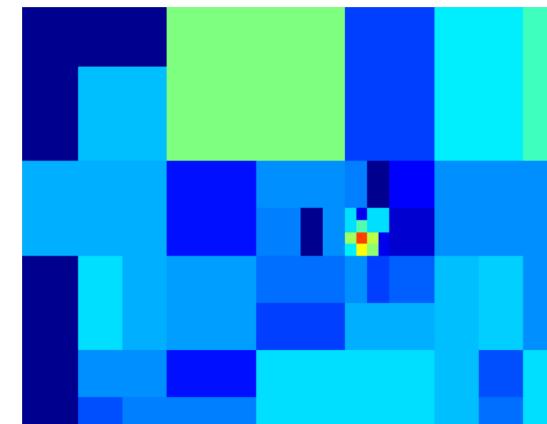
BB for $\arg \max_x S(x)$



Best-case: $O(P \log N)$



BB for $S(x) \geq -1$



Deformable Part Models (DPMs)

$$\mathbf{S}(\mathbf{x}) = \sum_{p=1}^P U_p(x_p) + B_p(x, x_p)$$

Local appearance

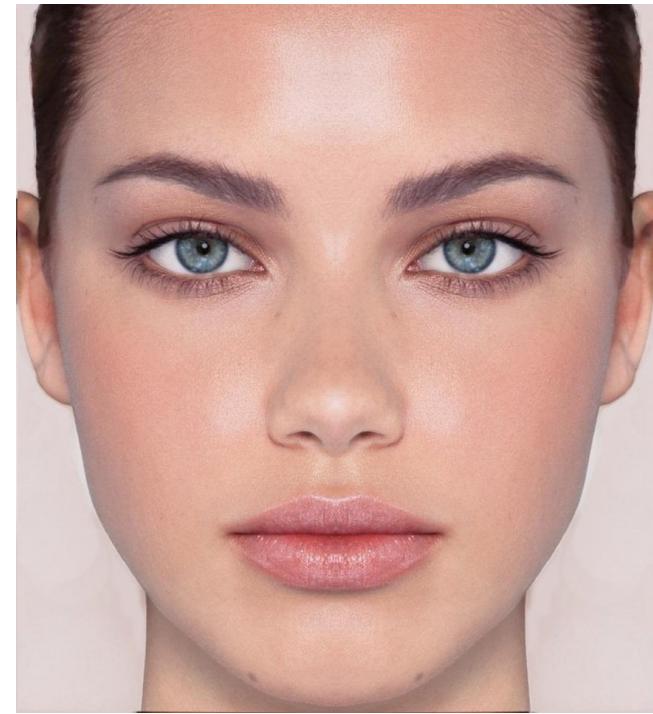
$$U_p(x_p) = \langle w_p, H(x_p) \rangle$$

Pairwise compatibility

$$B_p(x, x_p) = -(h - h_p - \hat{h}_p)^2 \eta - (v - v_p - \hat{v}_p)^2 \nu$$

Object score

$$S(x) = \sum_{p=1}^P \max_{x'} [U_p(x') + B_p(x, x')]$$



Bounding the DPM cost function

Wanted: $\max_{x \in X} S(x) \leq \bar{S}(X)$

$$S(x) = \sum_{p=1}^P m_p(x) \quad m_p(x) = \max_{x'} [U_p(x') + B_p(x, x')]$$

Property: $\max_{x \in X} h(x) + g(x) \leq \max_{x \in X} h(x) + \max_{x \in X} g(x)$

$$\max_{x \in X} \sum_{p=1}^P m_p(x) \leq \sum_{p=1}^P \max_{x \in X} m_p(x)$$



to-do

$$\max_{x \in X} S(x) \leq \sum_{p=1}^P \bar{m}_p(X) \doteq \bar{S}(X)$$

Bounding the source-to-domain contributions

$$\mu_d^s \doteq \max_{x \in X_d} \max_{x' \in X_s} U(x') + B(x', x)$$

Upper bound: $\mu_d^s \leq \max_{x' \in X_s} U(x') + \max_{x \in X_d} \max_{x' \in X_s} B(x', x)$

Intuitively: take best source point, put it at best possible source location

Quantities involved: $\bar{u}^s \doteq \max_{x' \in X_s} U(x')$, $\bar{g}_d^s \doteq \max_{x \in X_d} \max_{x' \in X_s} B(x', x)$

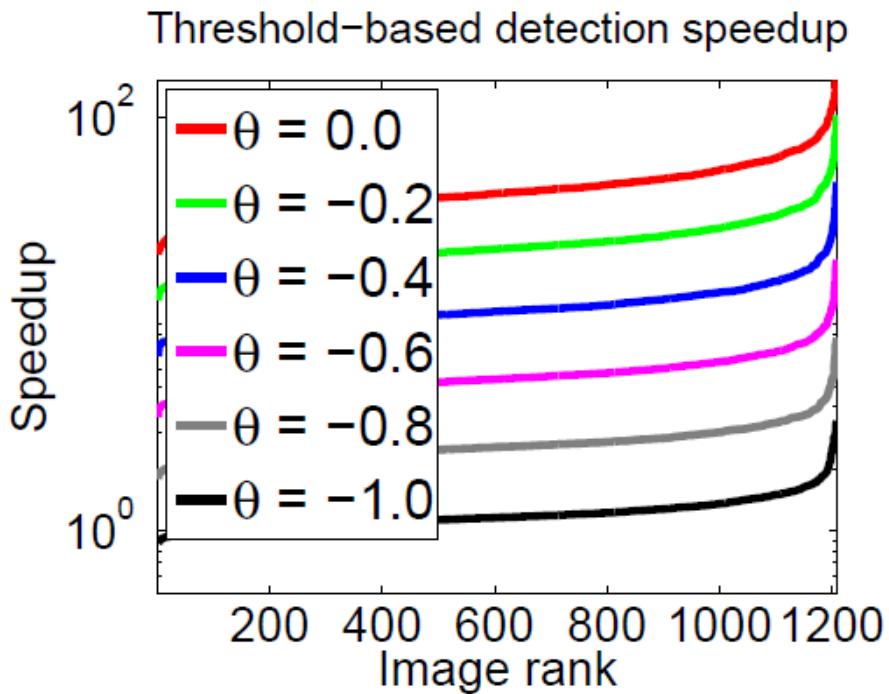
**Fine-to-coarse -
O(N)**

**Analytic
expression**

Pruning over combinations of s and d: Dual Trees

A.G. Gray and A.W. Moore. Nonparametric density estimation: Toward computational tractability. ICDM 2003

Results on Pascal VOC



1200 images, 20 Categories per image

Dual Trees

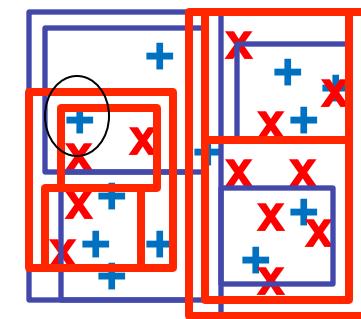
Data structure developed for Kernel Density estimation

$$P(x_j) = \sum_{i=1}^N w_i K(x_j, x_i), \quad x_i \in X_S, \quad x_j \in X_D \quad j = 1, \dots, M$$

Naïve solution: $O(|X_S||X_D|)$

Smarter: KD-tree for source points

Smartest: do this in batch mode for domain points



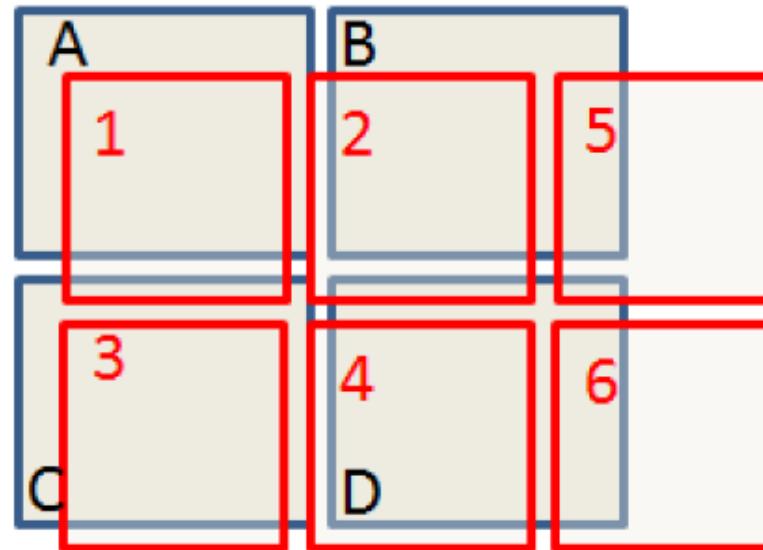
A.G. Gray and A.W. Moore. Nonparametric density estimation: Toward computational tractability. ICDM 2003

A. T. Ihler, E. B. Sudderth, W. T. Freeman, and A. S. Willsky. Efficient multiscale sampling from products of gaussian mixtures (NBP). In *NIPS*, 2003.

Breaking up the message bound

$$m(X) \doteq \max_{x \in X} \max_{x' \in X'} U(x') + B(x', x)$$

Partition domains: $X = \cup_{d \in D} X_d$, $X' = \cup_{s \in S} X_s$



Define: $\mu_d^s \doteq \max_{x \in X_d} \max_{x' \in X_s} U(x') + B(x', x)$

$$m(X) = \max_d \max_s \mu_d^s$$

Bounding the source-to-domain contributions

$$\mu_d^s \doteq \max_{x \in X_d} \max_{x' \in X_s} U(x') + B(x', x)$$

Constrained optimization: $\max_{x \in X_d} \left[\max_{x' \in X_s} U(x') + \max_{x'' \in X_s} B(x'', x) \right]$

s.t. $x'' = x'$

Removing the constraint gives upper bound:

$$\mu_d^s \leq \max_{x' \in X_s} U(x') + \max_{x \in X_d} \max_{x' \in X_s} B(x', x)$$

Intuitively: take best source point, put it at best possible source location

Quantities involved: $\bar{u}^s \doteq \max_{x' \in X_s} U(x')$, $\bar{g}_d^s \doteq \max_{x \in X_d} \max_{x' \in X_s} B(x', x)$

Bounding the source-to-domain contributions-II

$$\lambda_d^s = \min_{x \in X_d} \max_{x' \in X_s} U(x') + B(x', x)$$

Lower bound -i $\underline{\lambda}_{d,1}^s = \max_{x' \in X_s} U(x') + \min_{x \in X_d} \min_{x' \in X_s} B(x', x) \leq \lambda_d^s$

Intuitively: take best source point, put it at worst possible location

Lower bound -ii $\underline{\lambda}_{d,2}^s = \min_{x' \in X'_s} U(x') + \min_{x \in X_d} \max_{x' \in X_s} B(x', x) \leq \lambda_d^s$

Intuitively: take worst source point, put it at best possible location

$$\underline{\lambda}_d^s = \max(\underline{\lambda}_{d,1}^s, \underline{\lambda}_{d,2}^s)$$

Analytical expressions for geometric bounds

Pairwise term: $\mathcal{G}_{x,x'} = -H(h - h')^2 - V(v - v')^2$

$$\begin{aligned}\mathcal{G}_{\bar{d},\bar{s}} &\stackrel{\cdot}{=} \max_{x \in X_d} \max_{x' \in X_s} \mathcal{G}_{x,x'} \\ &= \max_{h \in X_d^h} \max_{h' \in X_s^h} -H(h - h')^2 + \max_{v \in X_d^v} \max_{v' \in X_s^v} -H(v - v')^2 \\ &= -Hh_{\underline{d},\underline{s}}^2 - Vv_{\underline{d},\underline{s}}^2, \quad \text{where}\end{aligned}$$

$$h_{\underline{d},\underline{s}} \stackrel{\cdot}{=} \min_{h \in X_d^h} \min_{h' \in X_s^h} |h - h'|, \quad v_{\underline{d},\underline{s}} \stackrel{\cdot}{=} \min_{v \in X_d^v} \min_{v' \in X_s^v} |v - v'|$$

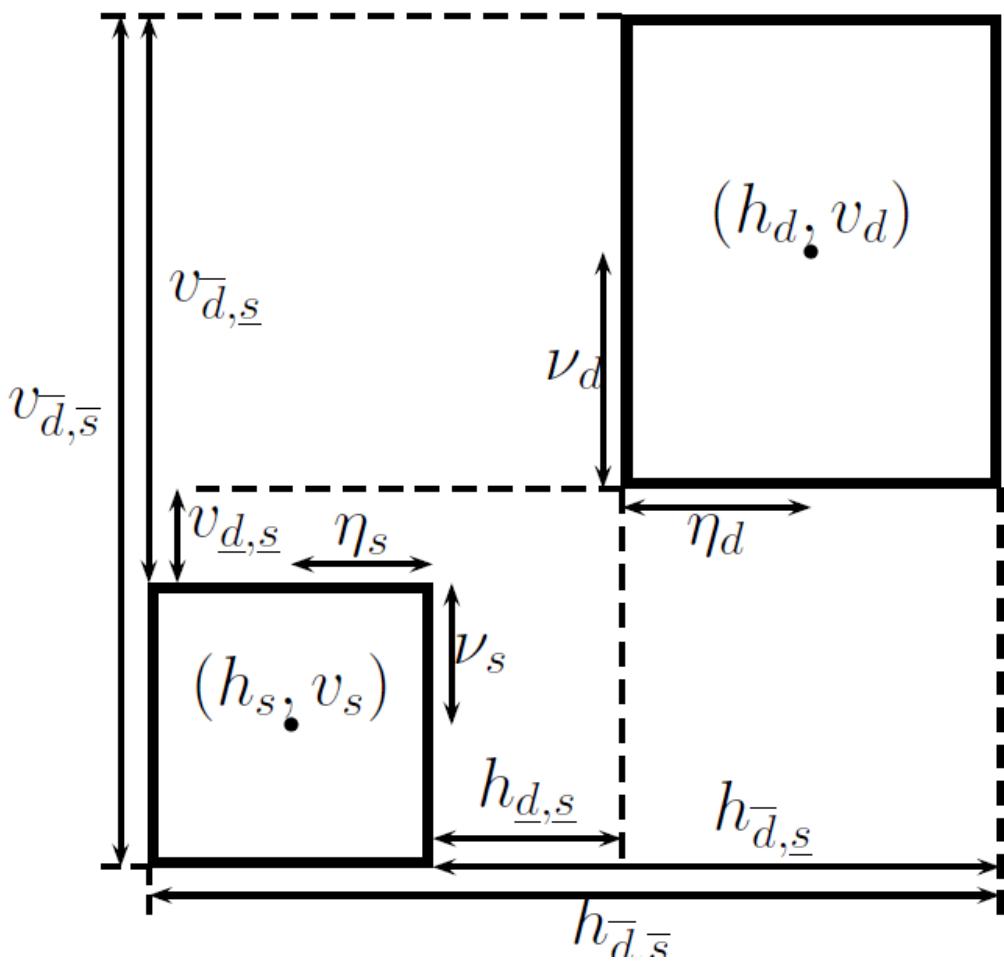
$$\mathcal{G}_{\underline{d},\bar{s}} \stackrel{\cdot}{=} \min_{x \in X_d} \max_{x' \in X_s} \mathcal{G}_{x,x'} = -Hh_{\bar{d},\underline{s}}^2 - Vv_{\bar{d},\underline{s}}^2, \quad \text{where}$$

$$h_{\bar{d},\underline{s}} \stackrel{\cdot}{=} \max_{h \in X_d^h} \min_{h' \in X_s^h} |h - h'|, \quad v_{\bar{d},\underline{s}} \stackrel{\cdot}{=} \max_{v \in X_d^v} \min_{v' \in X_s^v} |v - v'|$$

$$\mathcal{G}_{\underline{d},\bar{s}} \stackrel{\cdot}{=} \min_{x \in X_d} \min_{x' \in X_s} \mathcal{G}_{x,x'} = -Hh_{\bar{d},\bar{s}}^2 - Vv_{\bar{d},\bar{s}}^2, \quad \text{where}$$

$$h_{\bar{d},\bar{s}} \stackrel{\cdot}{=} \max_{h \in X_d^h} \max_{h' \in X_s^h} |h - h'|, \quad v_{\bar{d},\bar{s}} \stackrel{\cdot}{=} \max_{v \in X_d^v} \max_{v' \in X_s^v} |v - v'|$$

Exploit rectangular domain



$$\begin{aligned}
 h_{\bar{d},\bar{s}} &= (h_d + \eta_d) - (h_s + \eta_s) \\
 h_{\bar{d},s} &= (h_d - \eta_d) - (h_s + \eta_s) \\
 h_{d,\underline{s}} &= (h_d + \eta_d) - (h_s - \eta_s) \\
 v_{\bar{d},\bar{s}} &= (v_d + \nu_d) - (v_s + \nu_s) \\
 v_{\bar{d},s} &= (v_d - \nu_d) - (v_s + \nu_s) \\
 v_{d,\underline{s}} &= (v_d + \nu_d) - (v_s - \nu_s)
 \end{aligned}$$

Done so far

$$m(X) \doteq \max_{x \in X} \max_{x' \in X'} U(x') + B(x', x)$$

$$X = \cup_{d \in D} X_d, \quad X' = \cup_{s \in S} X_s$$

$$\mu_d^s \doteq \max_{x \in X_d} \max_{x' \in X_s} U(x') + B(x', x)$$

$$m(X) = \max_d \max_s \mu_d^s$$

$$\overline{m}(X) = \max_d \max_s \overline{\mu}_d^s$$

Large intervals: loose bounds

Small intervals: tight bounds, but maximization over multiple terms

$$|S| = \frac{|X|}{|X_s|}$$

Need to prune

Dual Recursion

Source-/ Domain-tree: KD-trees for source/domain points, respectively.

Goal: keep manageable the computation of $\overline{m}(X) = \max_d \max_s \overline{\mu}_d^s$

Key idea: entertain list of s that can possibly contribute to any d

$$\mathcal{S}_d = \{s_i\}$$

Originally: d is Domain-tree root, its supporter is Source-tree root

$$\mathcal{S}_0 = \{s_0\}$$

When splitting (branching) d, split also its supporters

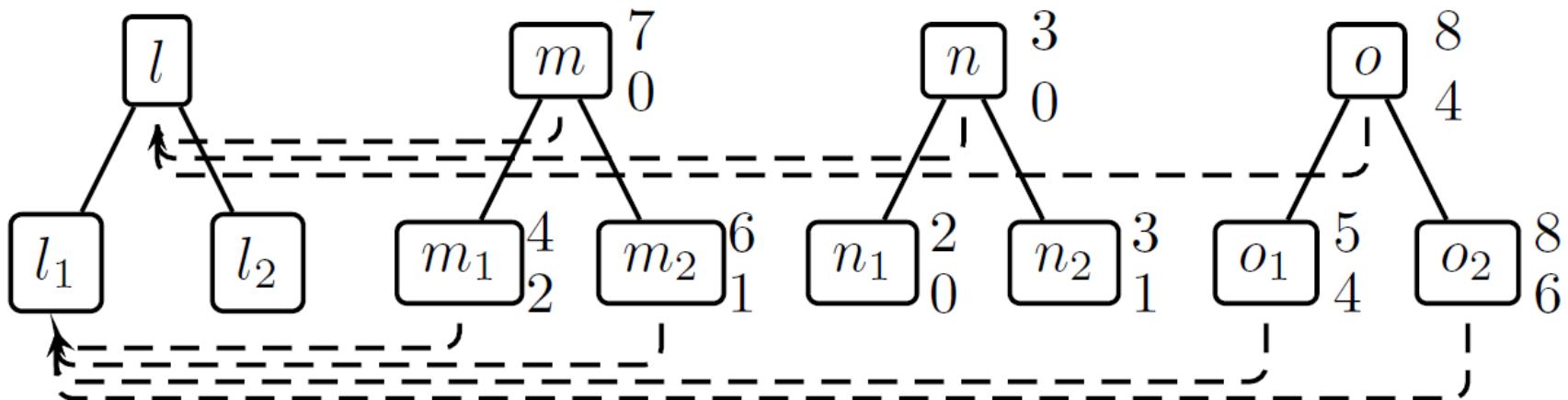
Prune supporters based on upper and lower bounds

Supporter pruning

$\mathcal{S}_l = \{m, n, o\}$ Source nodes m; n; o ‘support’ domain node l

Recurse: $\mathcal{S}_{ch(l)} = \text{Prune} [ch(m) \cup ch(n) \cup ch(o)]$

Pruning criterion: $\overline{\mu}_d^l < \max_{j \in \mathcal{S}_d} \underline{\lambda}_d^j$



DTBB demonstration



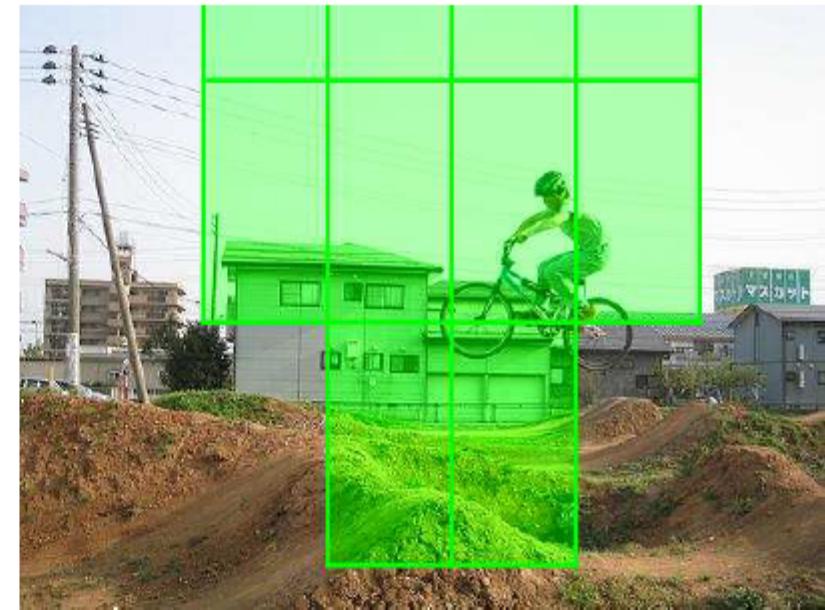
DTBB demonstration



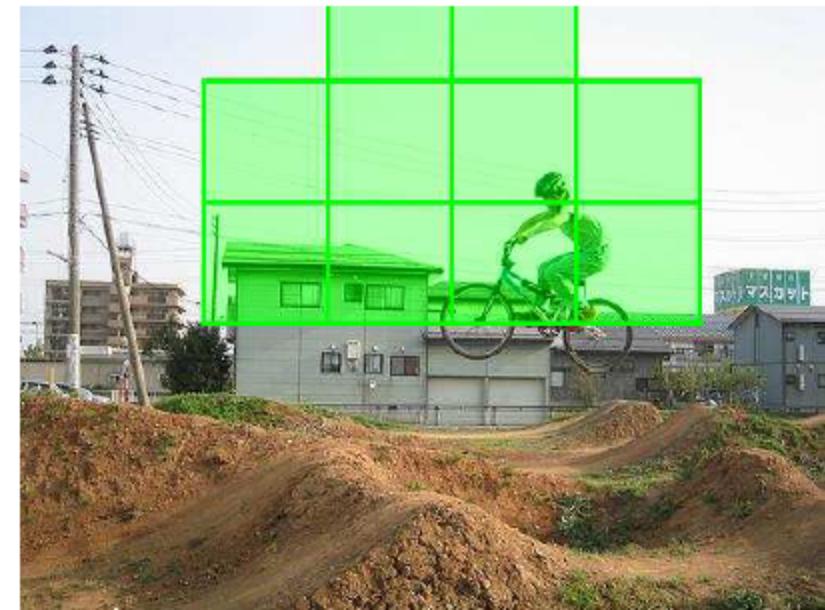
DTBB demonstration



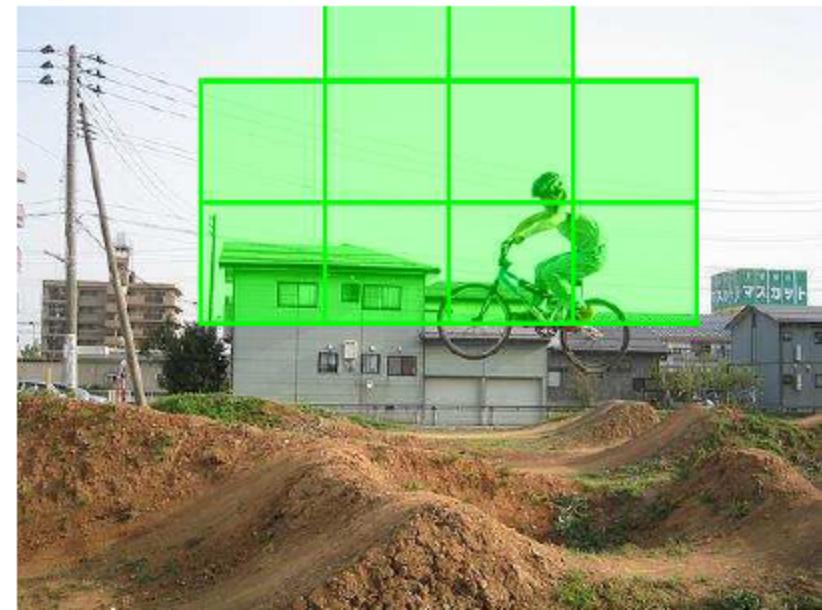
DTBB demonstration



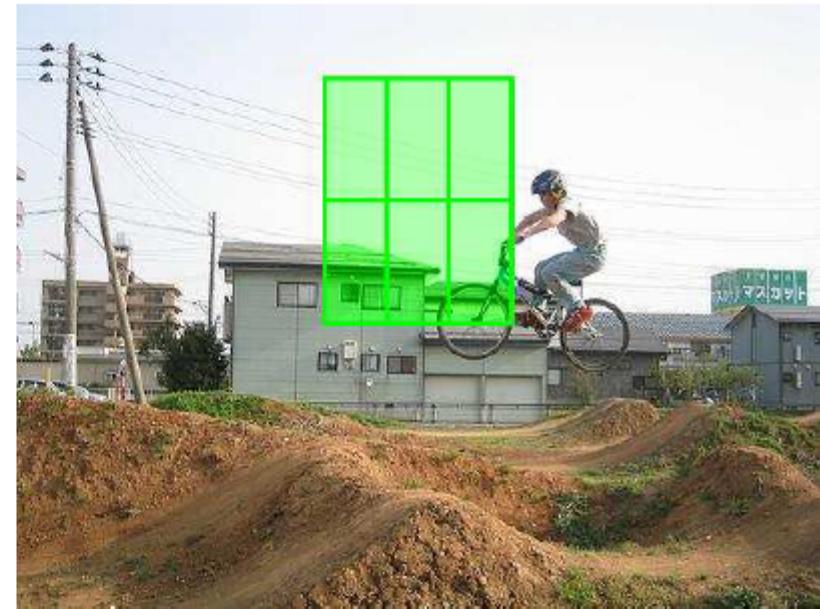
DTBB demonstration



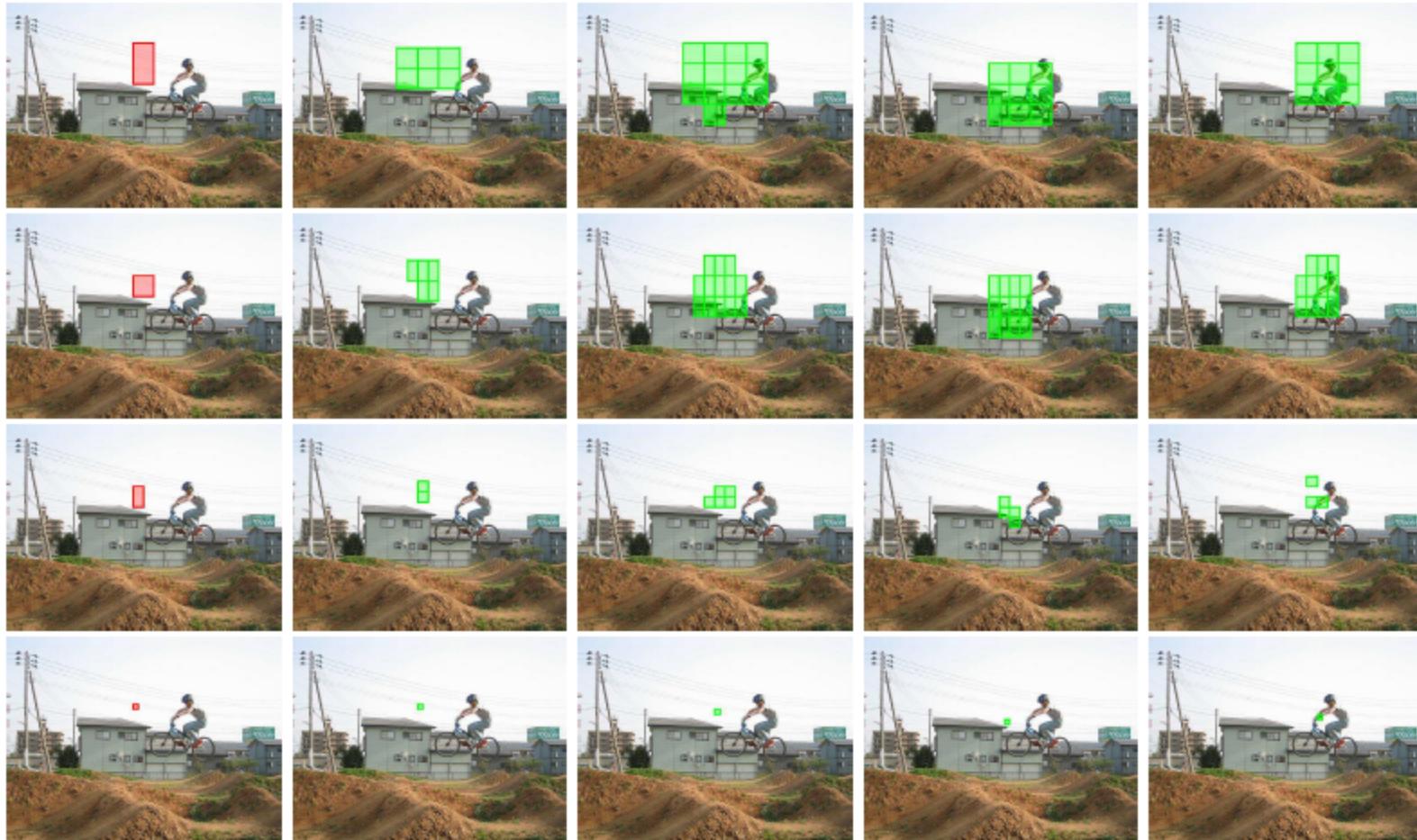
DTBB demonstration



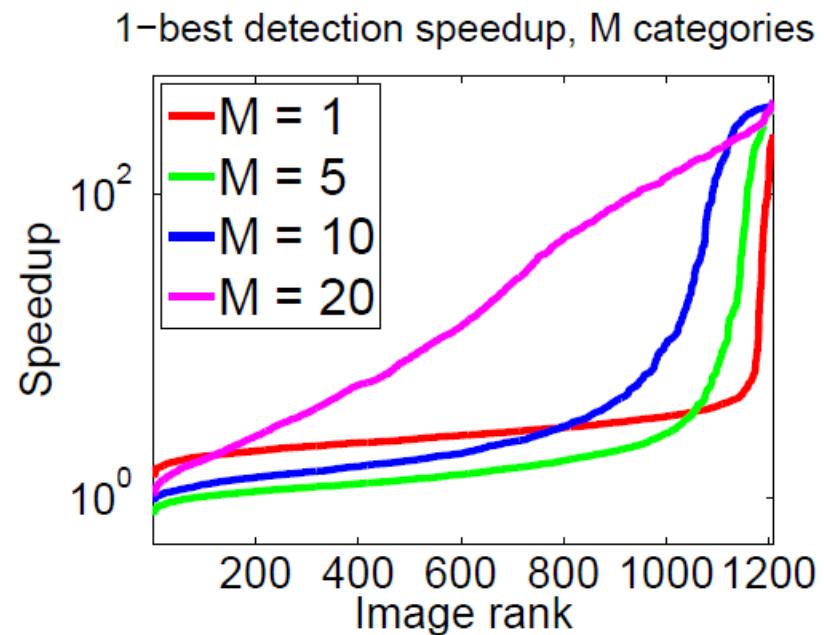
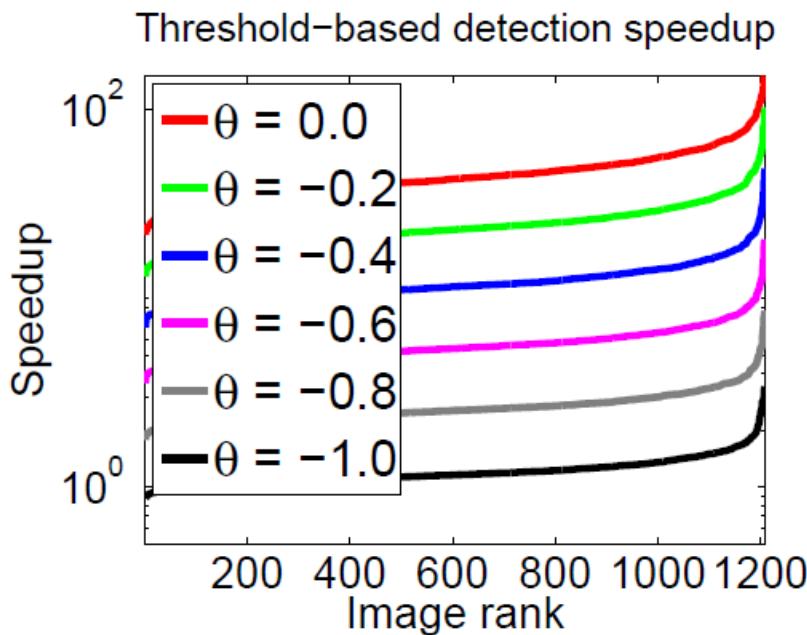
DTBB demonstration



Supporter pruning demonstration



Results on Pascal VOC



1200 images, 20 Categories per image

Results on Pascal VOC

	Our algorithm	[4]
Unary terms	13.20 ± 1.49	159.41 ± 15.82
KD-trees	1.72 ± 0.21	0.00 ± 0.00
Detection, $\theta = 0.0$	0.25 ± 0.07	10.74 ± 1.02
Detection, $\theta = -.2$	0.47 ± 0.12	10.74 ± 1.02
Detection, $\theta = -.4$	0.93 ± 0.22	10.74 ± 1.02
Detection, $\theta = -.6$	1.95 ± 0.42	10.74 ± 1.02
Detection, $\theta = -.8$	4.17 ± 0.84	10.74 ± 1.02
Detection, $\theta = -1$	9.14 ± 1.79	10.74 ± 1.02
Detection, 1-best	0.41 ± 0.08	10.74 ± 1.02
Detection, 5-best	0.47 ± 0.09	10.74 ± 1.02
Detection, 10-best	0.48 ± 0.10	10.74 ± 1.02

1200 images, 20 Categories per image

Object detection with Deformable Part Models (DPMs)

$$U_p(x') = \langle \mathbf{w}_p, \mathbf{H}(x') \rangle$$



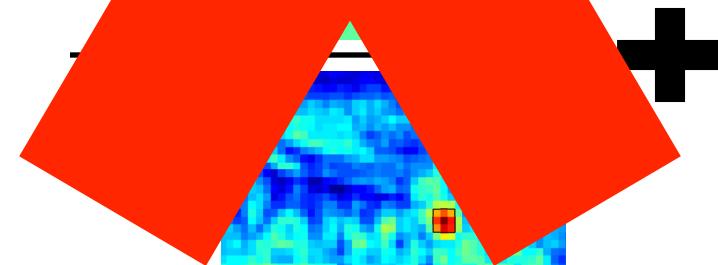
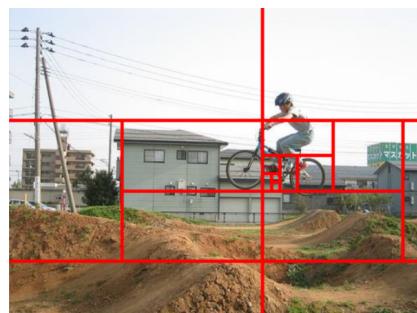
$$\max_{x'} [U_p(x') + B_p(x, x')]$$

$$p = 1$$

⋮

DTBB, NIPS 11

$$p = P$$



$$S(x) = \sum_{p=1}^P \max_{x'} [U_p(x') + B_p(x, x')]$$

Accelerating detection with DPMs

Efficient detection with DPMs

P F Felzenszwalb R B Girshick and D A McAllester Cascade object detection with DPMs CVPR 2010

B. Sapp, A. Toshev, and B. Taskar. Cascaded models for articulated pose estimation, ECCV, 2010

M. Pedersoli, A. Vedaldi, and J. Gonzalez. A coarse-to-fine approach for object detection, CVPR 2011

I. Kokkinos, Rapid DPM Detection using Dual-Tree Branch-and-Bound, NIPS 2011

Efficient part score computation

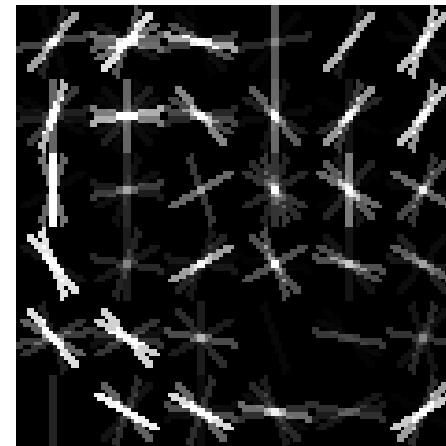
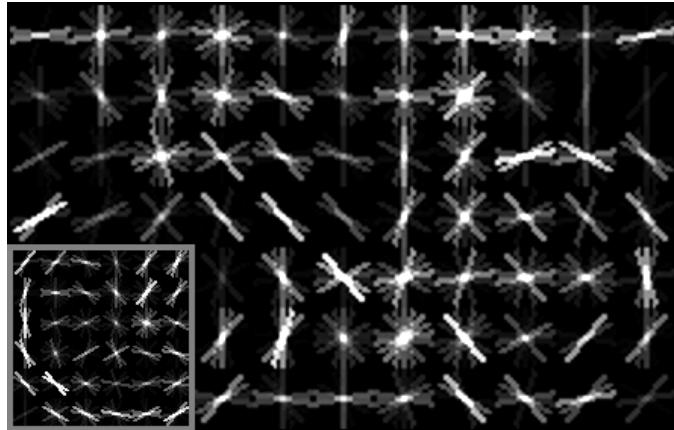
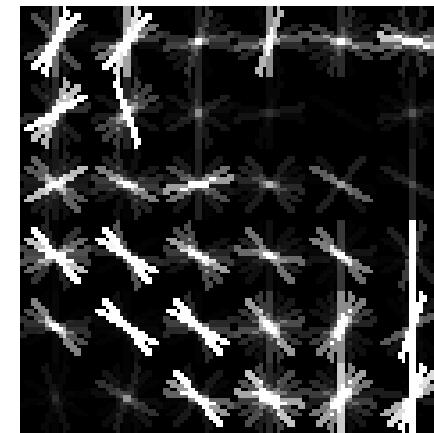
H. Pirsiavash and D. Ramanan, Steerable Part Models, CVPR 2012

A. Vedaldi and A. Zisserman, Sparse Kernel Maps and Faster Product Quantization Learning, CVPR 2012

H.O. Song, S. Zickler, T. Althoff, R. Girschick, M. Fritz, C. Geyer, P. Felzenszwalb, T. Darrell, Sparselet Models for Efficient Multiclass Object Detection, ECCV 2012

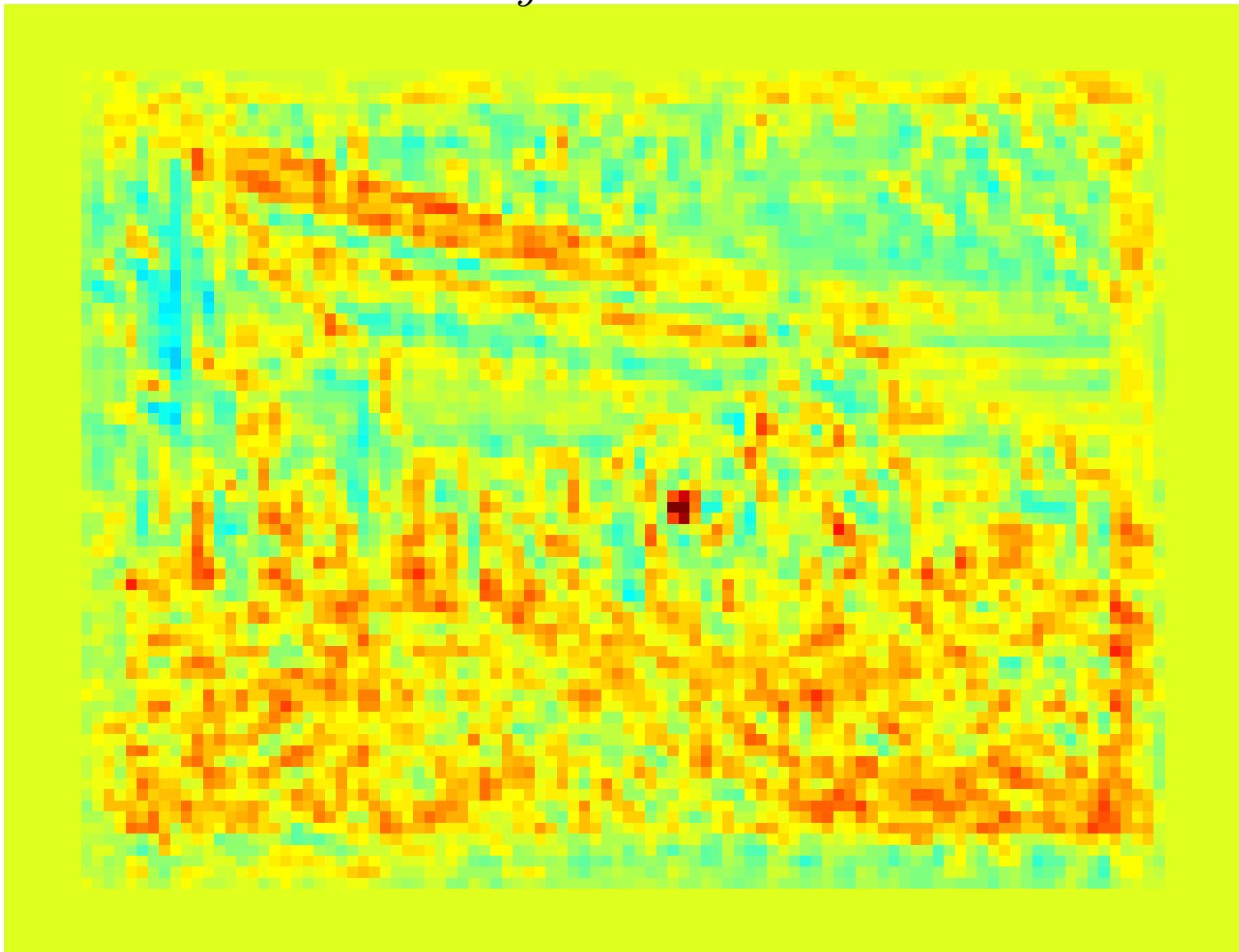
C. Dubout and F. Fleuret. Exact Acceleration of Linear Object Detectors, ECCV 2012

Part score computation

 $\mathbf{w}[y]$  $\mathbf{h}[x + y]$

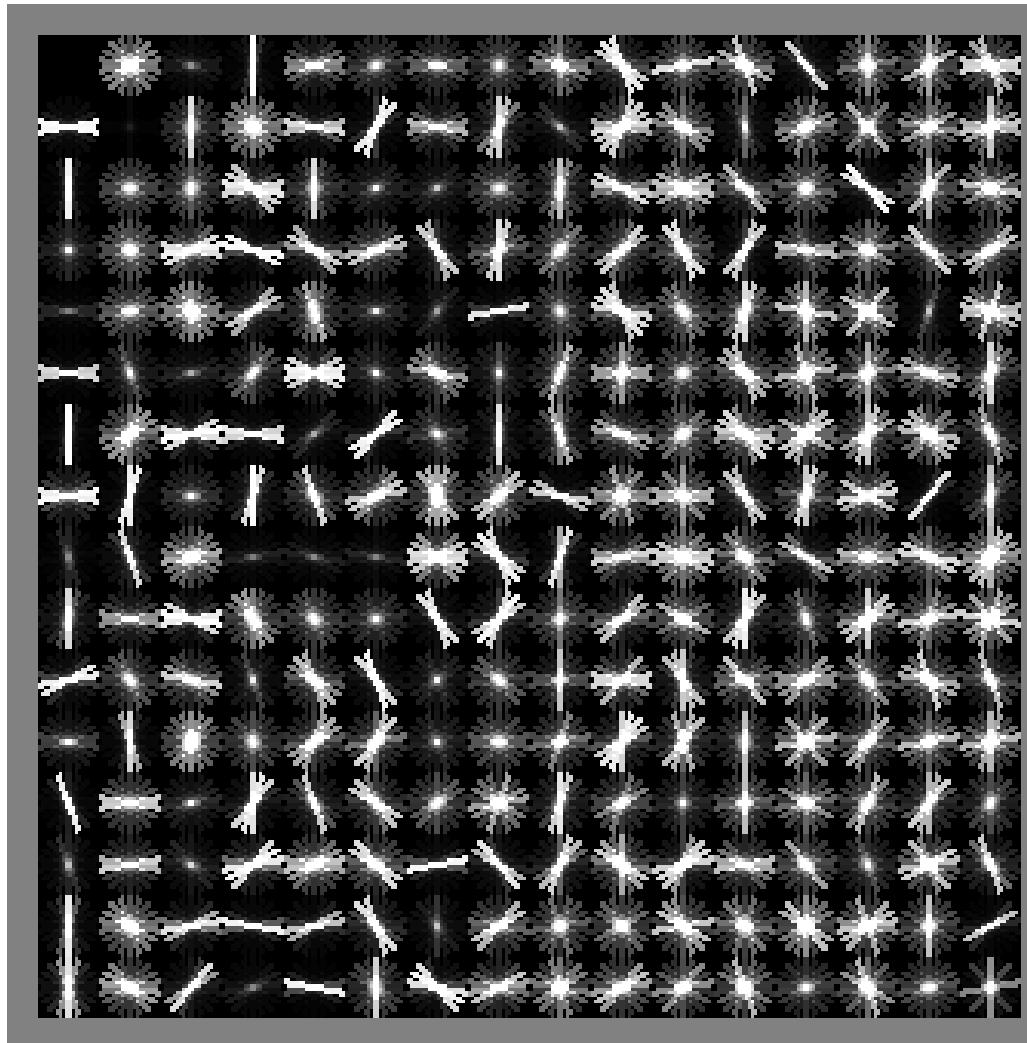
$$s[x] = \sum_y \langle \mathbf{h}[x + y], \mathbf{w}[y] \rangle$$

Part scores $s[x] = \sum_y \langle \mathbf{h}[x + y], \mathbf{w}[y] \rangle$



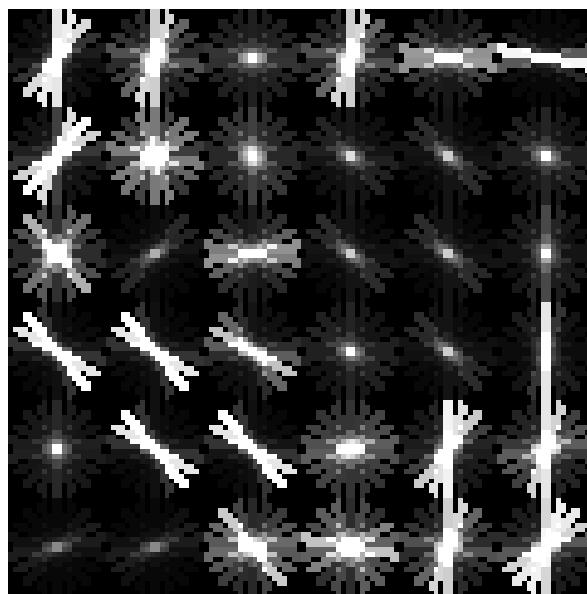
HOG cell quantization: visual ‘letters’

$$\mathcal{C} = \{C_1, \dots, C_{256}\}$$



HOG feature quantization

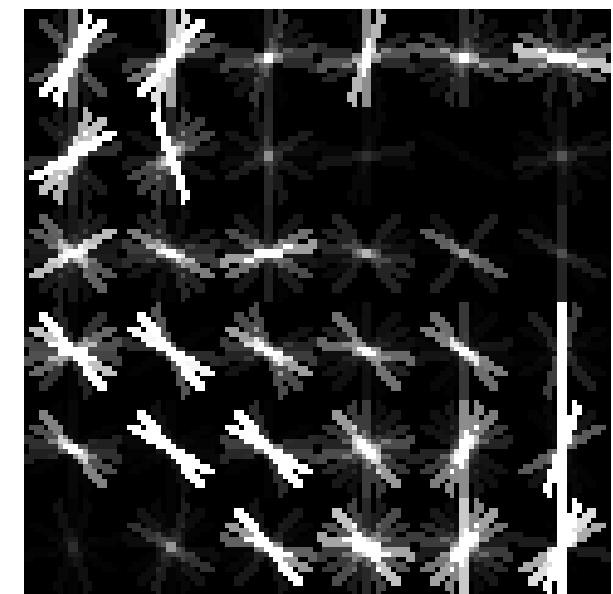
HOG detail



Codebook indices

60	199	39	199	25	62
121	143	132	45	129	85
209	70	64	129	129	117
210	210	200	85	129	118
3	210	210	20	185	115
63	63	186	242	199	155

Quantized HOG



$$\mathbf{h}[x]$$

$$i[x] = \arg \min_k d(\mathbf{h}[x], C_k)$$

$$\hat{\mathbf{h}}[x] = C_{i[x]}$$

Efficient inner product approximation

$$\left\langle \begin{array}{c} \text{[Image of a textured surface]} \\ , \end{array} \right\rangle \approx \left\langle \begin{array}{c} \text{[Image of a textured surface]} \\ , \end{array} \right\rangle$$

$$s[x] \simeq \hat{s}[x]$$

$$\sum_y \langle \mathbf{h}[x+y], \mathbf{w}[y] \rangle \simeq \sum_y \langle \hat{\mathbf{h}}[x+y], \mathbf{w}[y] \rangle$$

Efficient inner product approximation

$$\langle \begin{matrix} \text{image} \\ \text{image} \end{matrix}, \begin{matrix} \text{image} \\ \text{image} \end{matrix} \rangle \simeq \langle \begin{matrix} \text{image} \\ \text{image} \end{matrix}, \begin{matrix} \text{image} \\ \text{image} \end{matrix} \rangle$$

$$\langle \mathbf{h}[x+y], \mathbf{w}[y] \rangle \simeq \langle \hat{\mathbf{h}}[x+y], \mathbf{w}[y] \rangle$$

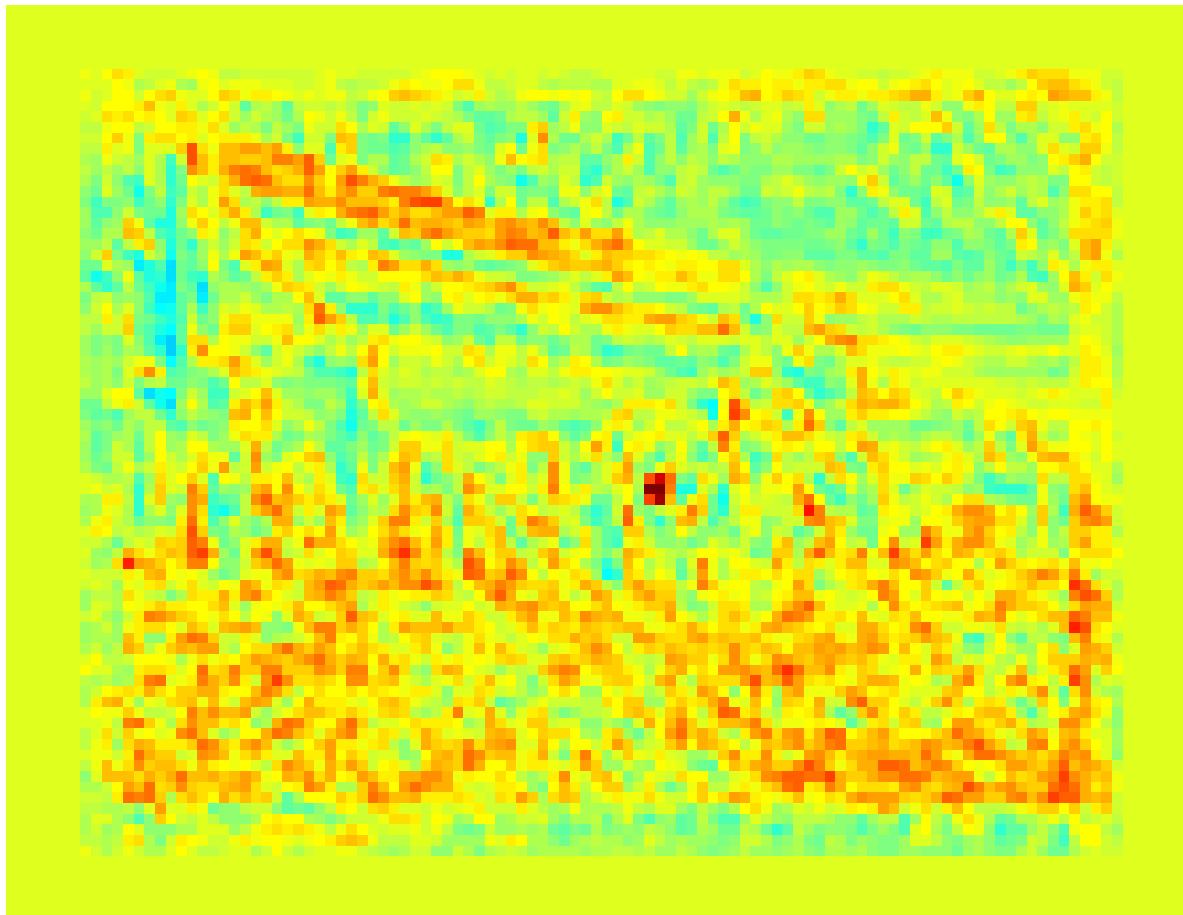
$$= \langle C_{I[x+y]}, \mathbf{w}[y] \rangle$$

$$= \Pi[I[x+y], y]$$

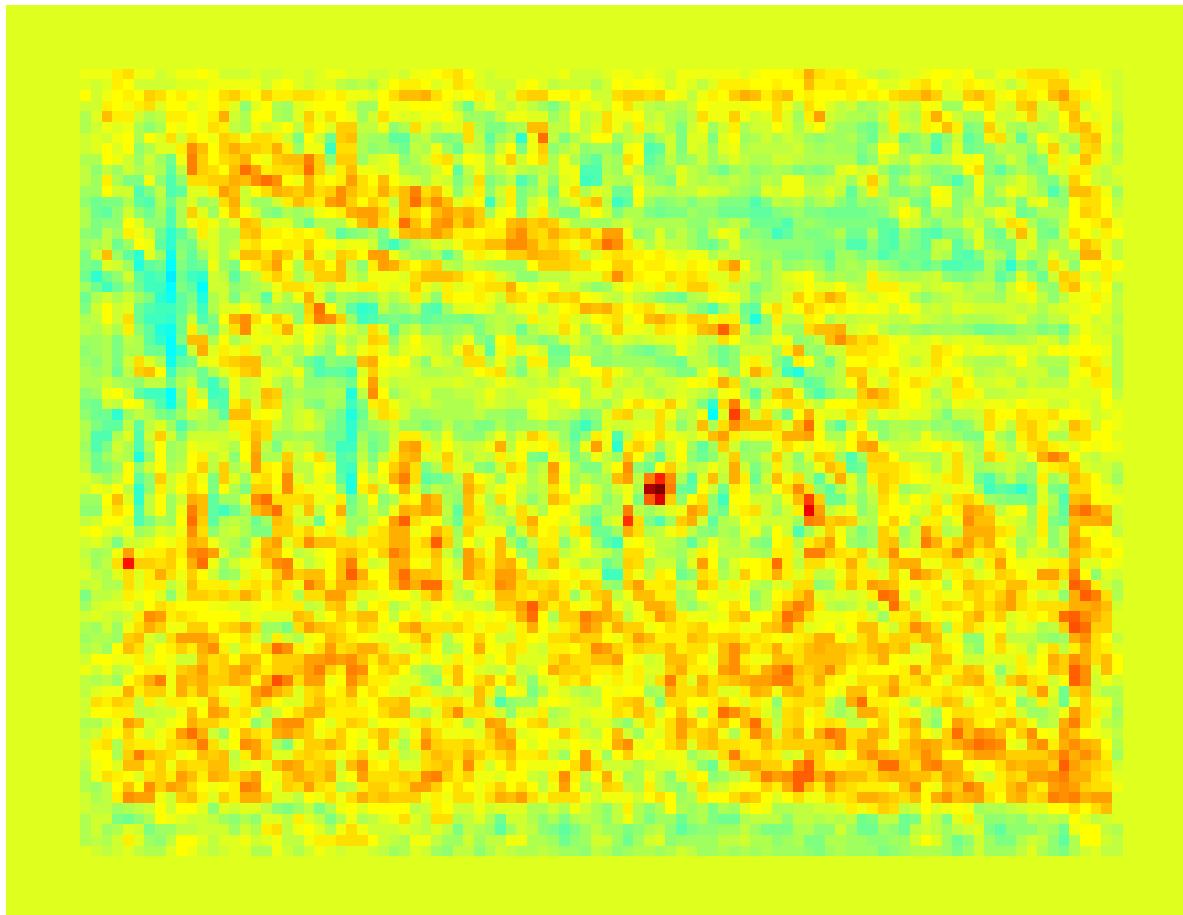
$$\Pi[k, y] = \langle C_k, \mathbf{w}_y \rangle$$

$$s[x] \simeq \hat{s}[x] = \sum_y \Pi[I[x+y], y]$$

Lookup-based estimate demonstration: $s[x]$



Lookup-based estimate demonstration: $\hat{s}[x]$



Part-level approximation error

$$\begin{aligned}\epsilon \doteq s - \hat{s} &= \left\langle \begin{array}{c} \text{[Image of a scene]} \\ \text{[Image of a scene]} \end{array} - \begin{array}{c} \text{[Image of a scene]} \\ \text{[Image of a scene]} \end{array}, \begin{array}{c} \text{[Image of a scene]} \\ \text{[Image of a scene]} \end{array} \right\rangle \\ &= \sum_y \underbrace{\langle \mathbf{h}[y] - \hat{\mathbf{h}}[y], \mathbf{w}[y] \rangle}_{e[y]}\end{aligned}$$

Cell-level approximation error

$$\begin{aligned} e[y] &= \langle \mathbf{h}[y] - \hat{\mathbf{h}}[y], \mathbf{w}[y] \rangle = \langle \text{[image]} - \text{[image]}, \text{[image]} \rangle \\ &= \langle \mathbf{e}[y], \mathbf{w}[y] \rangle \\ &= \sum_{f=1}^{32} \mathbf{e}_y[f] \mathbf{w}_y[f] \end{aligned}$$

Chebyshev inequality

For any zero-mean random variable, and any value of α :

$$P(|X| > \alpha) \leq \frac{E\{X^2\}}{\alpha^2}$$

Equivalently, with probability of error smaller than p_e :

$$X \in \left[-\sqrt{\frac{E\{X^2\}}{p_e}}, \sqrt{\frac{E\{X^2\}}{p_e}} \right]$$

Chebyshev inequality-II

For a weighted sum of i.i.d. zero-mean random variables:

$$X' = \sum_{k=1}^K w_k X_k$$

with probability of error smaller than p_e :

$$X' \in \left[-\sqrt{\frac{(\sum_k w_k^2) E\{X^2\}}{p_e}}, \sqrt{\frac{(\sum_k w_k^2) E\{X^2\}}{p_e}} \right]$$

Chebyshev inequality for cell-level error

$$e[y] = \sum_{f=1}^F \mathbf{e}_y[f] \mathbf{w}_y[f]$$

with probability of error smaller than p_e :

$$e_y \in \left[-\sqrt{\frac{\|\mathbf{w}[y]\|^2 \|\mathbf{e}[y]\|^2}{p_e F}}, \sqrt{\frac{\|\mathbf{w}[y]\|^2 \|\mathbf{e}[y]\|^2}{p_e F}} \right]$$

Chebyshev inequality for part-level error

$$\epsilon = \hat{s} - s = \sum_y e[y]$$

with probability of error smaller than p_e :

$$\epsilon \in \left[-\sqrt{\frac{\sum_y \|\mathbf{w}[y]\|^2 \|\mathbf{e}[y]\|^2}{p_e F}}, \sqrt{\frac{\sum_y \|\mathbf{w}[y]\|^2 \|\mathbf{e}[y]\|^2}{p_e F}} \right]$$



with probability of error smaller than p_e :

$$s \in \left[\hat{s} - \sqrt{\frac{\sum_y \|\mathbf{w}[y]\|^2 \|\mathbf{e}[y]\|^2}{p_e F}}, \hat{s} + \sqrt{\frac{\sum_y \|\mathbf{w}[y]\|^2 \|\mathbf{e}[y]\|^2}{p_e F}} \right]$$

Recap

Lookup-based approximation:

$$s[x] \simeq \hat{s}[x] = \sum_y \Pi[I[x + y], y]$$

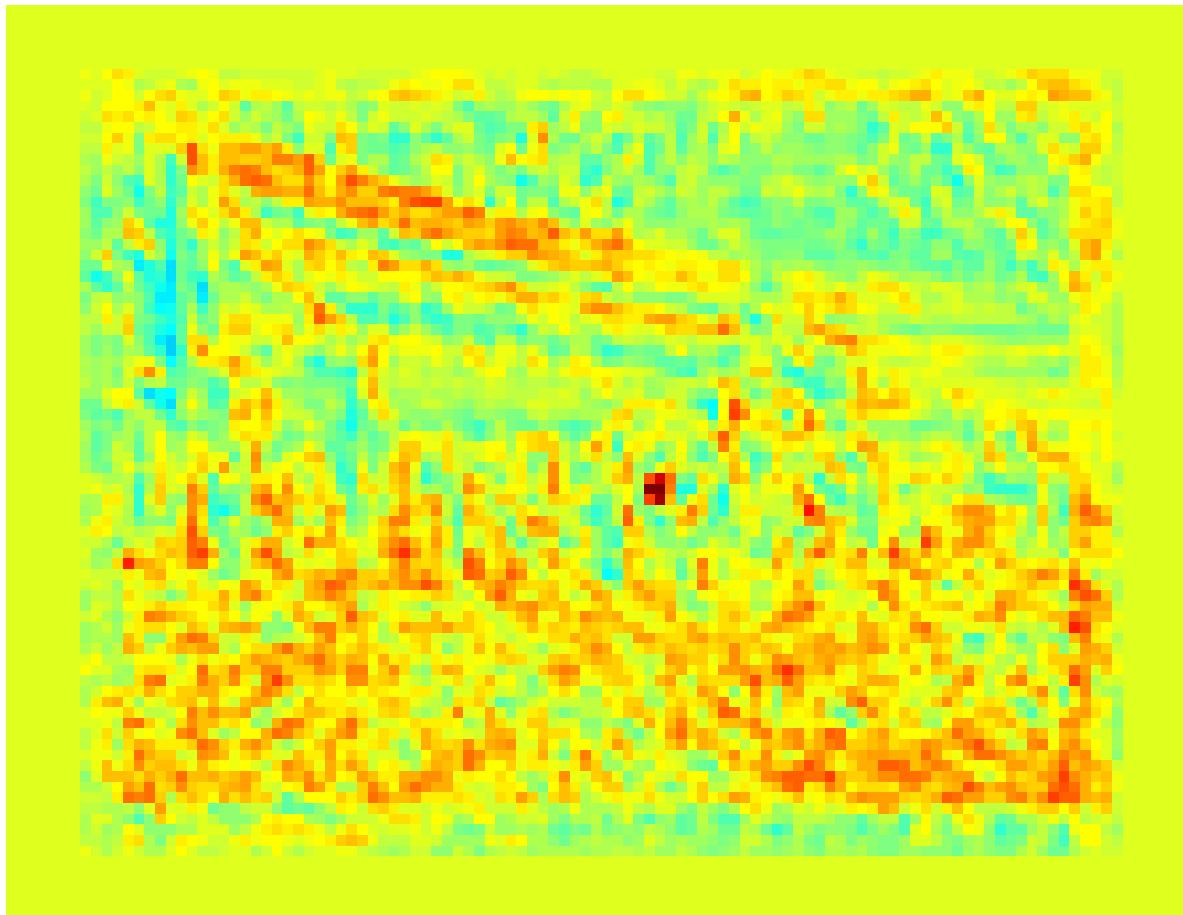
With probability of error at most p_e :

$$\underline{s}[x] \leq s[x] \leq \bar{s}[x]$$

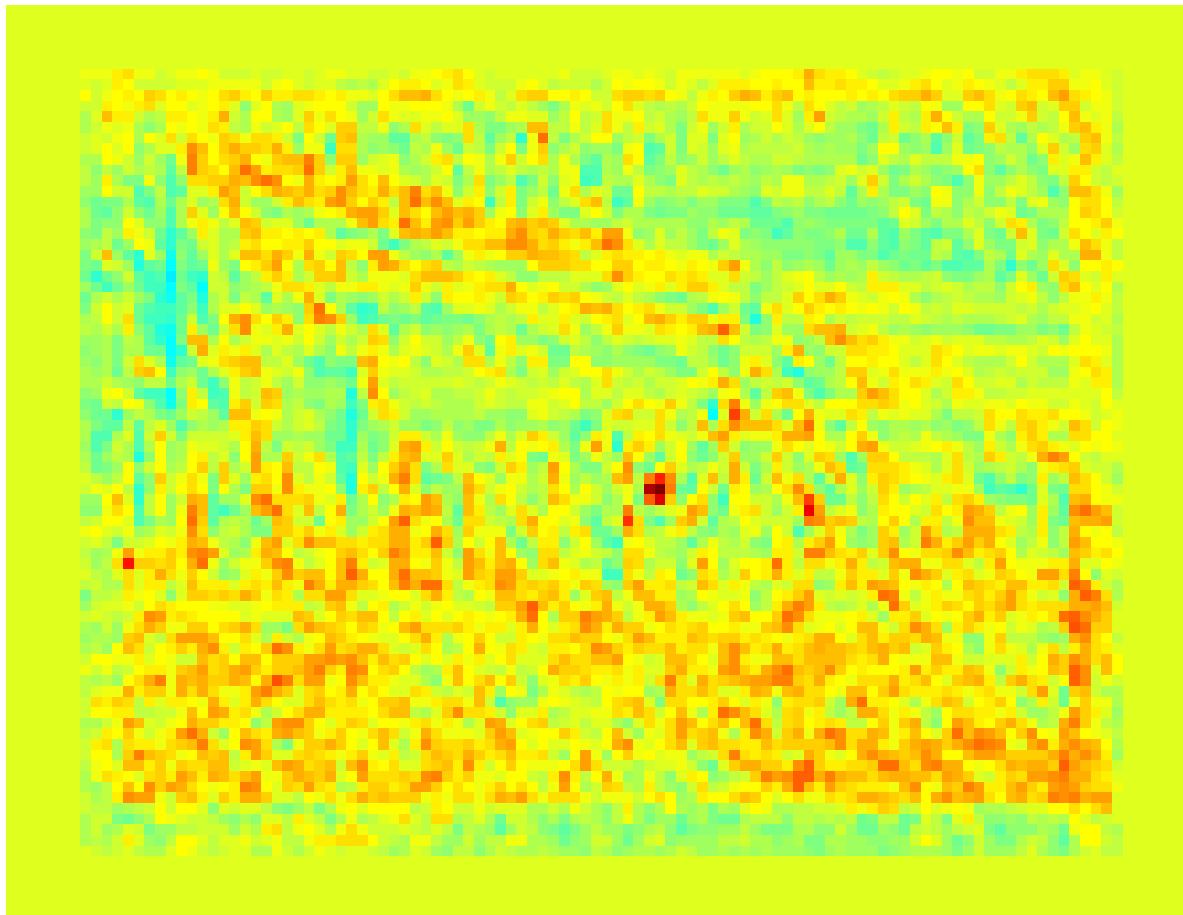
$$\underline{s}[x] = \hat{s}[x] - \sqrt{\frac{\sum_y \|\mathbf{w}[y]\|^2 \|\mathbf{e}[x + y]\|^2}{p_e F}}$$

$$\bar{s}[x] = \hat{s}[x] + \sqrt{\frac{\sum_y \|\mathbf{w}[y]\|^2 \|\mathbf{e}[x + y]\|^2}{p_e F}}$$

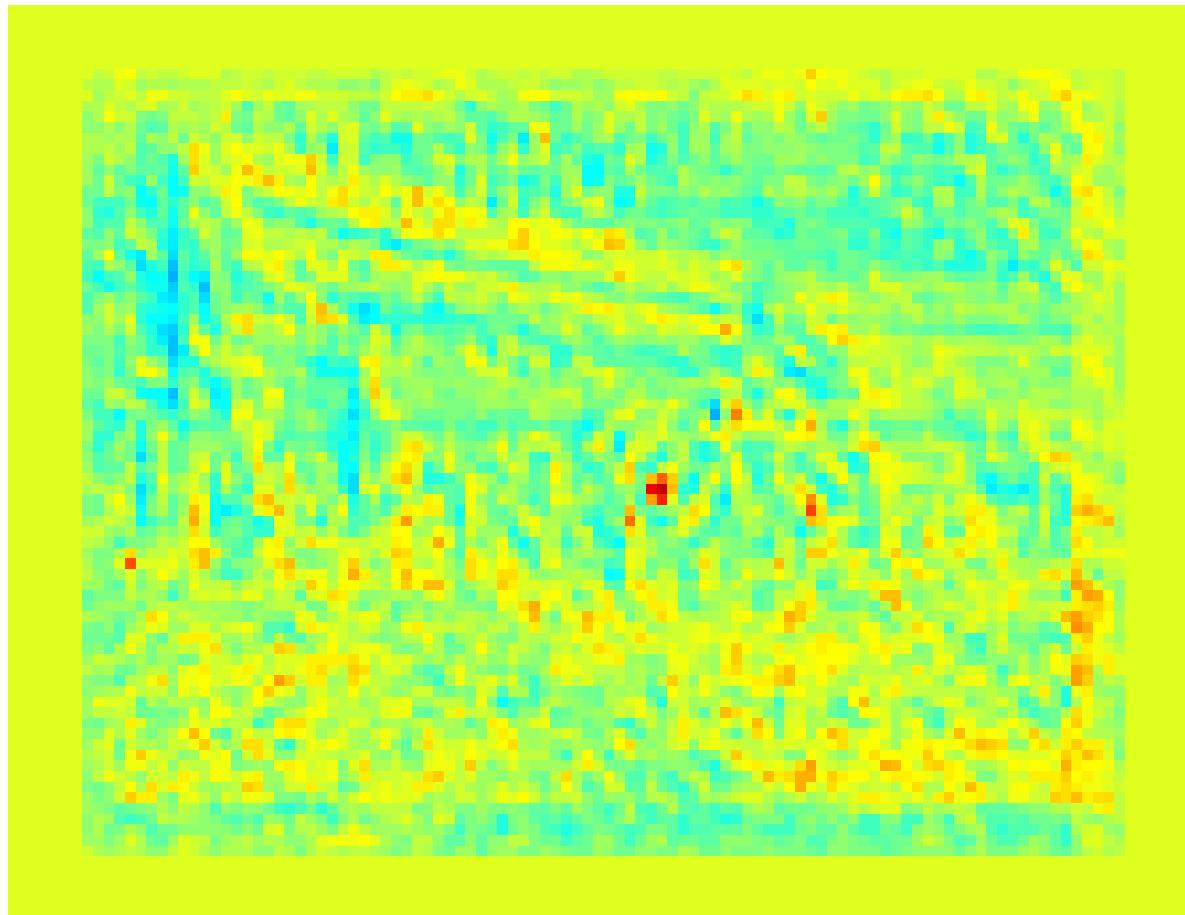
Bound demonstration: $s[x]$



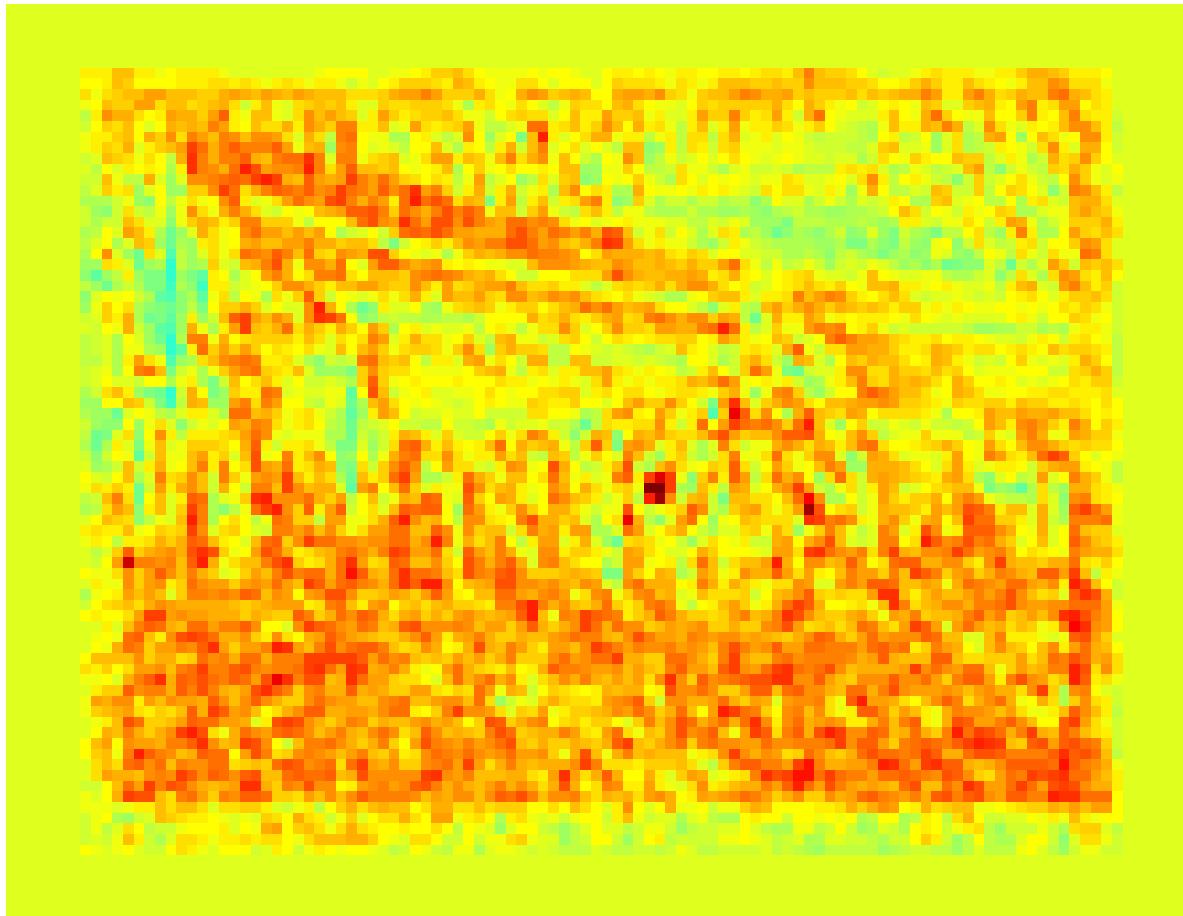
Bound demonstration: $\hat{s}[x]$



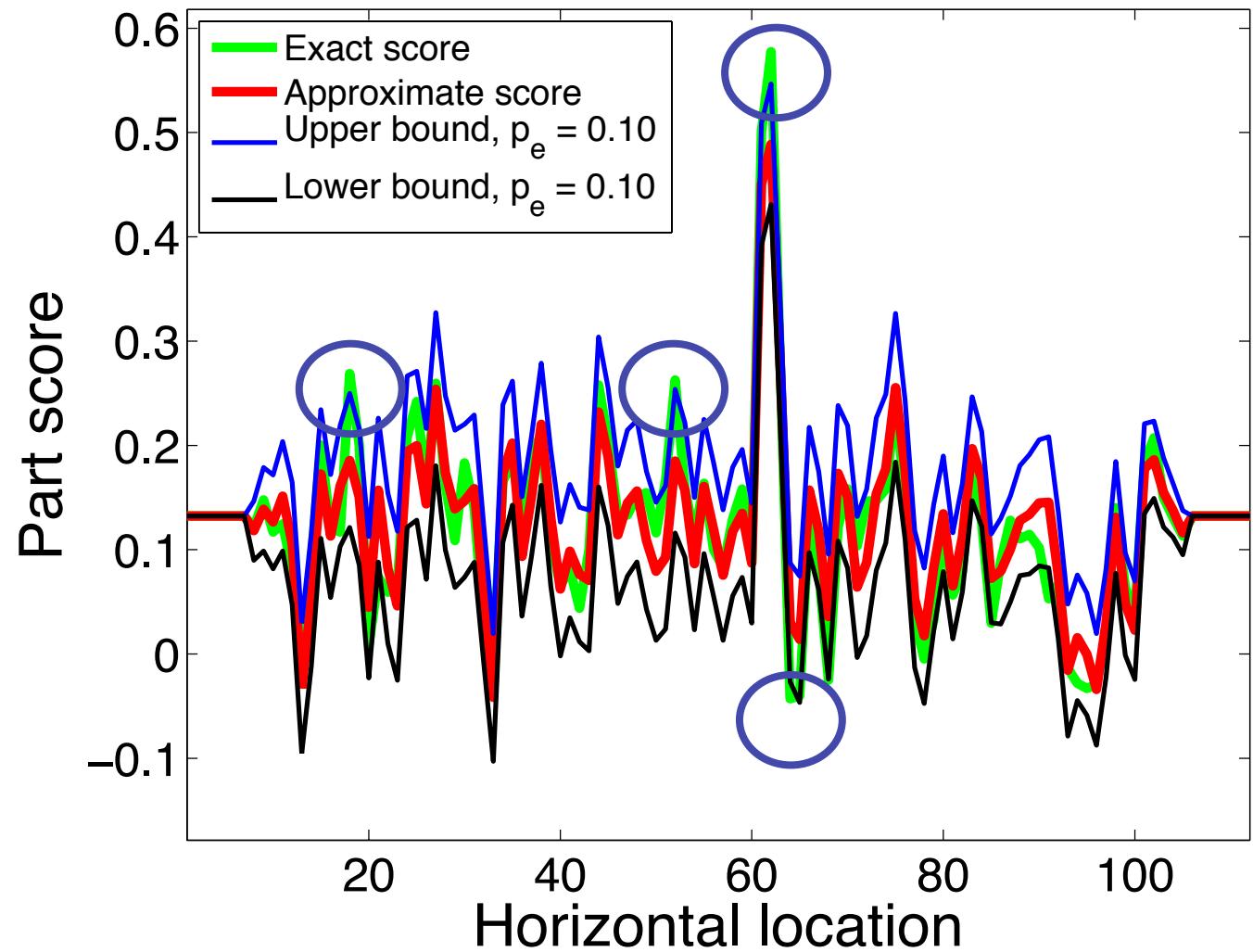
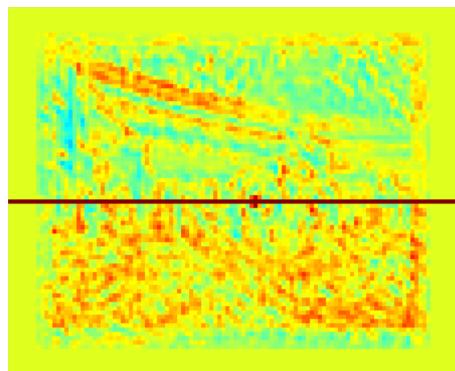
Bound demonstration: $\underline{s}[x]$, $p_e = .05$



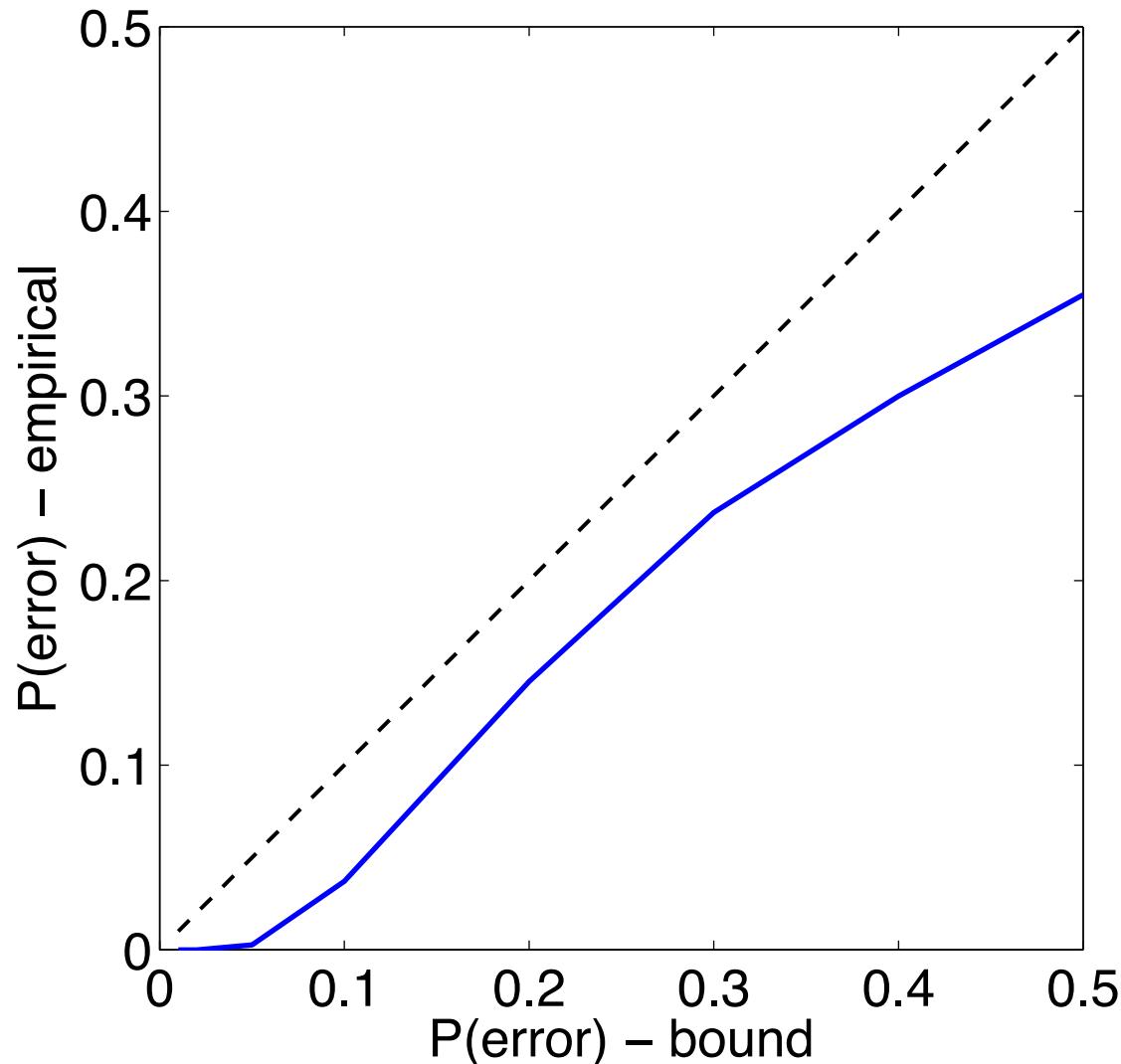
Bound demonstration: $\bar{s}[x]$, $p_e = .05$



Bound demonstration for varying confidence

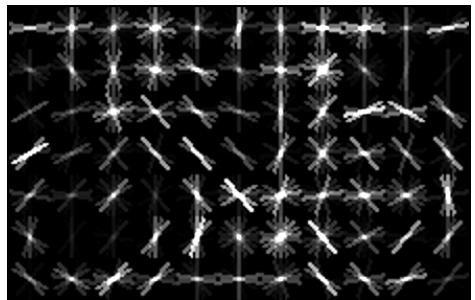


Bound tightness

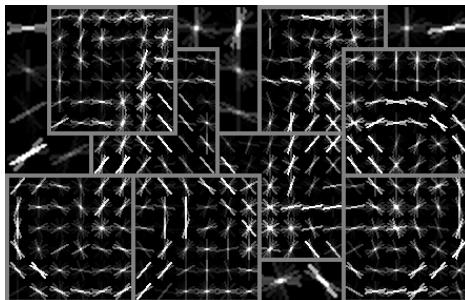


Accelerating detection with DPMs

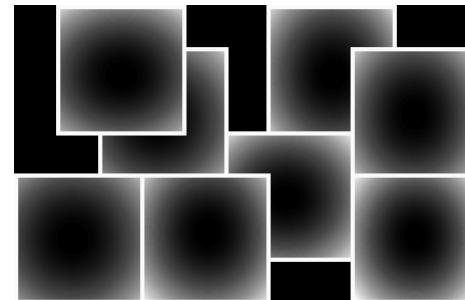
This work



w_1

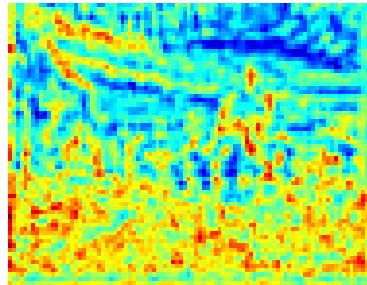


$w_2 \dots w_P$



$B_2 \dots B_P$

$$U_p(x) = \langle w_p, H(x) \rangle$$



⋮ GDT

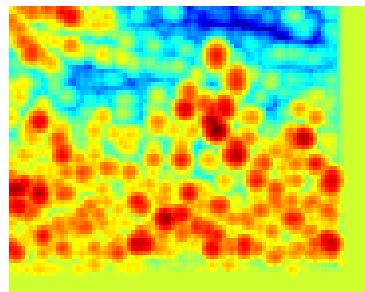
$$p = 3$$

⋮

$$p = 5$$

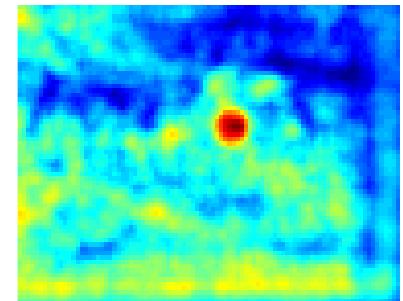
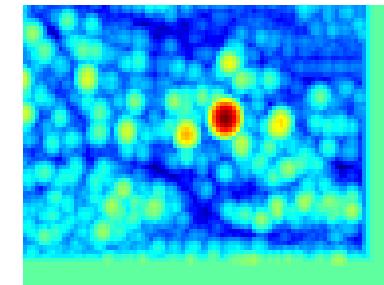
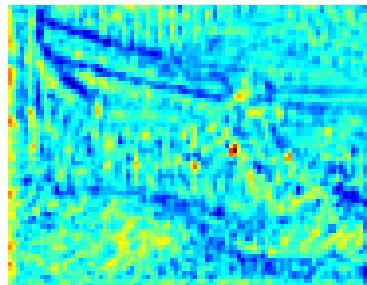
⋮

$$\mu_p(x) = \max_{x'} [U_p(x') + B_p(x, x')]$$



DTBB, NIPS 2011

$$S(x) = \sum_{p=1}^P \mu_p(x)$$

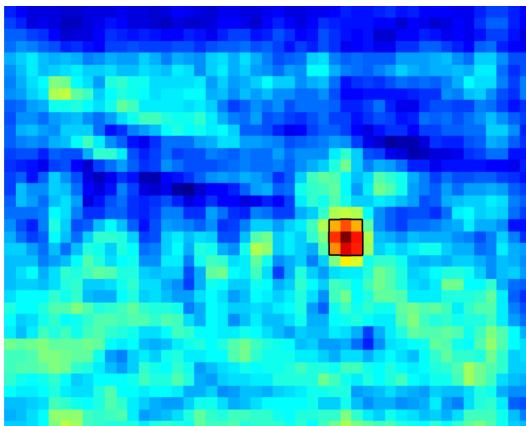


Dual-Tree Brand-and-Bound

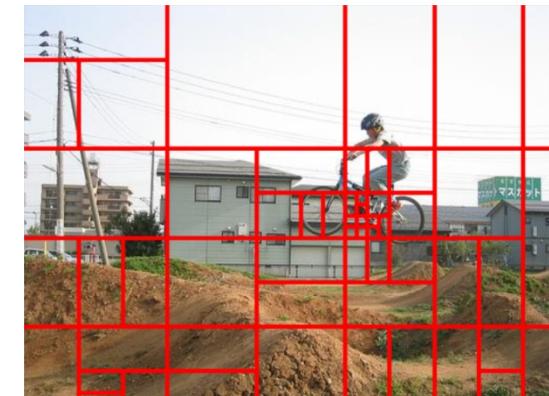
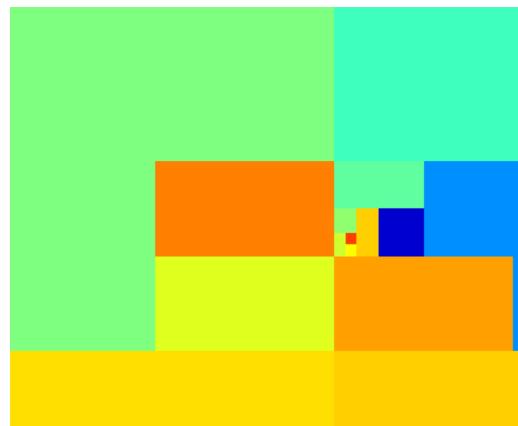
Input & Detection result



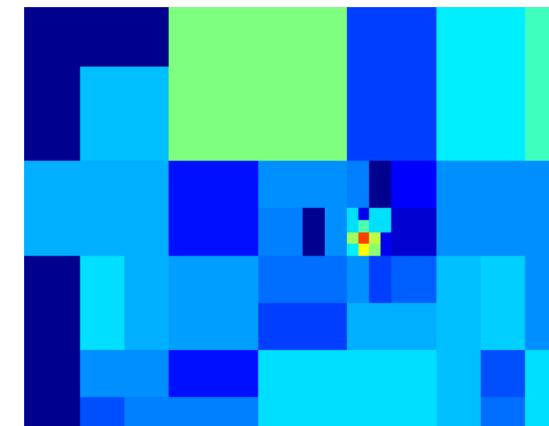
Detector score $S(x)$



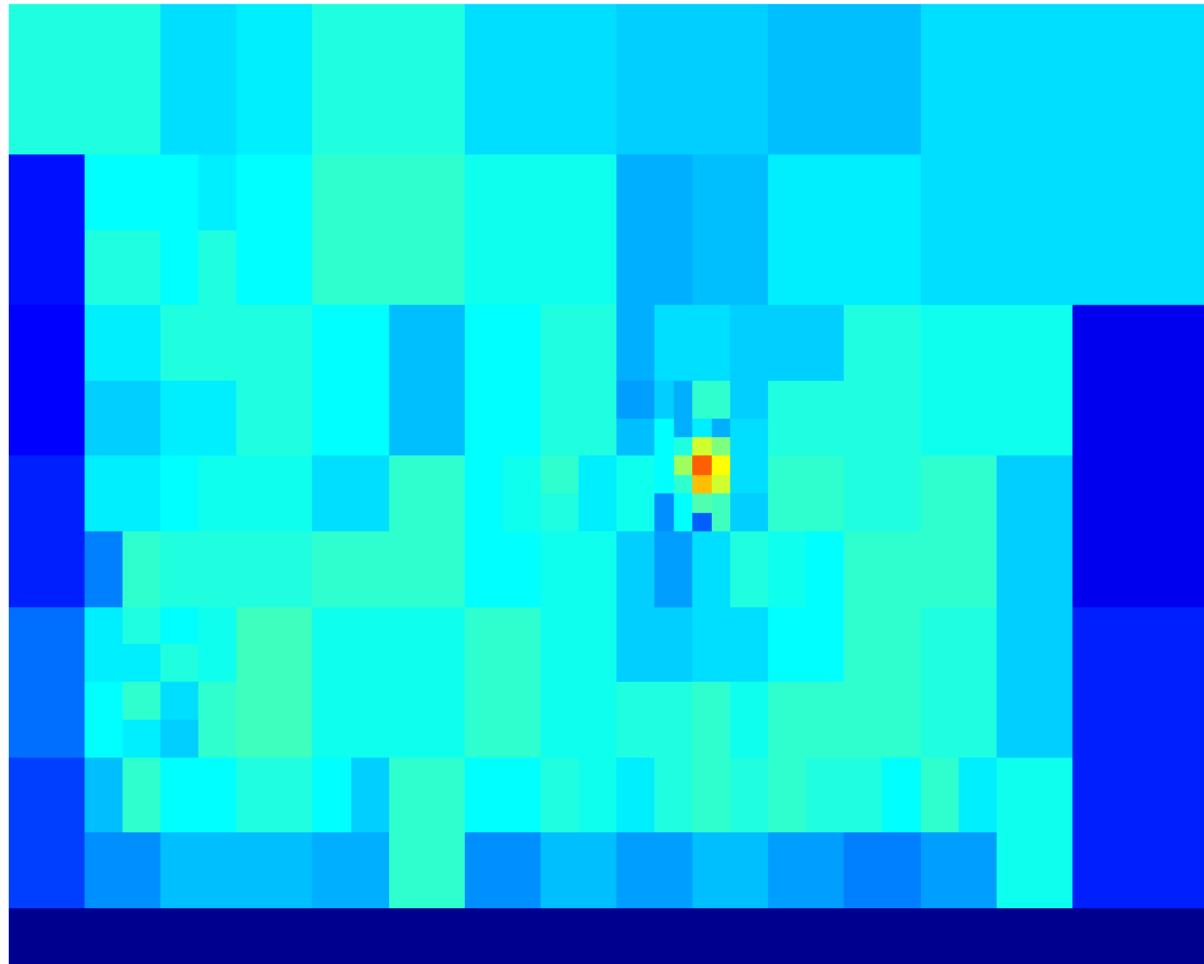
BB for $\arg \max_x S(x)$



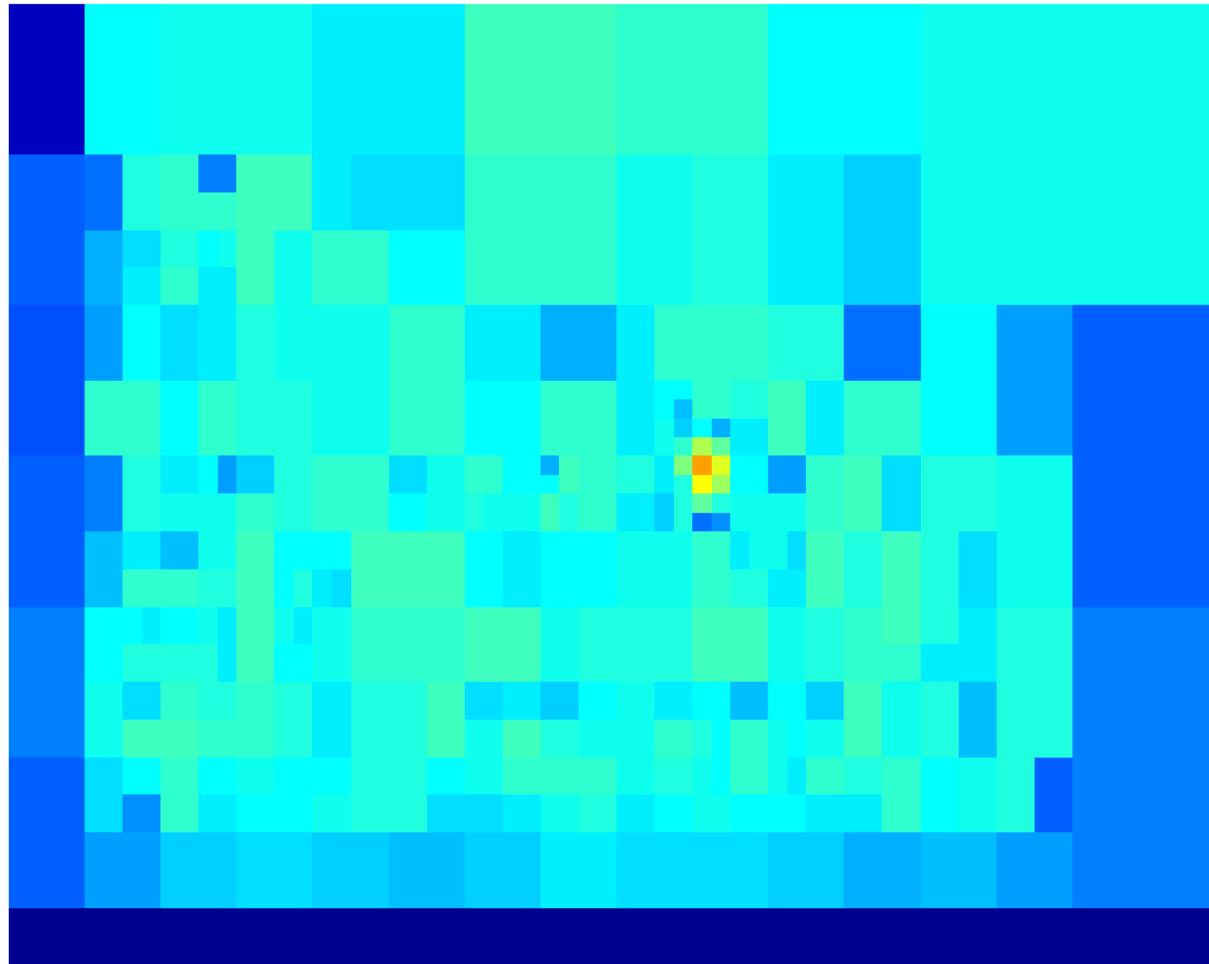
BB for $S(x) \geq -1$



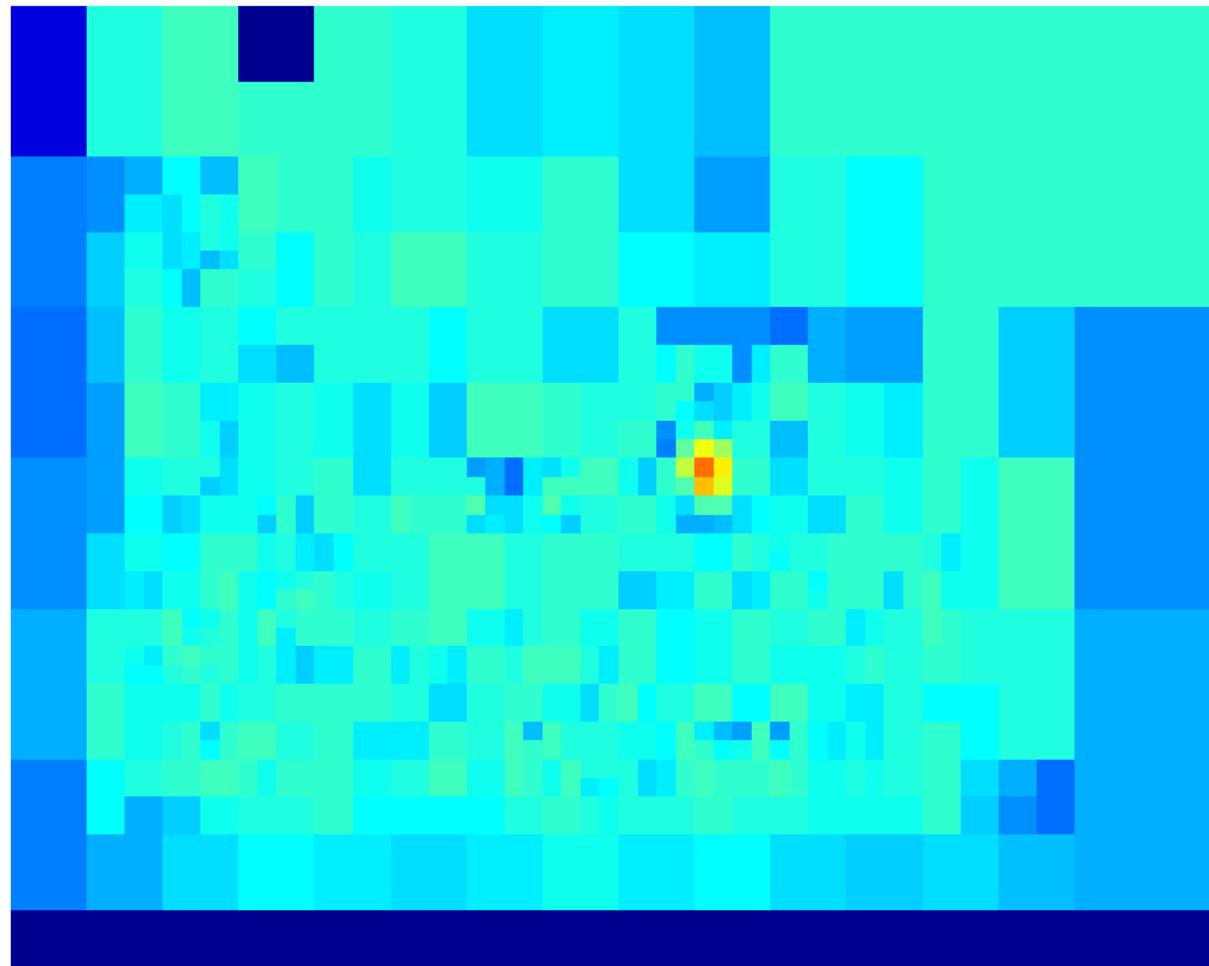
DTBB results: exact part scores



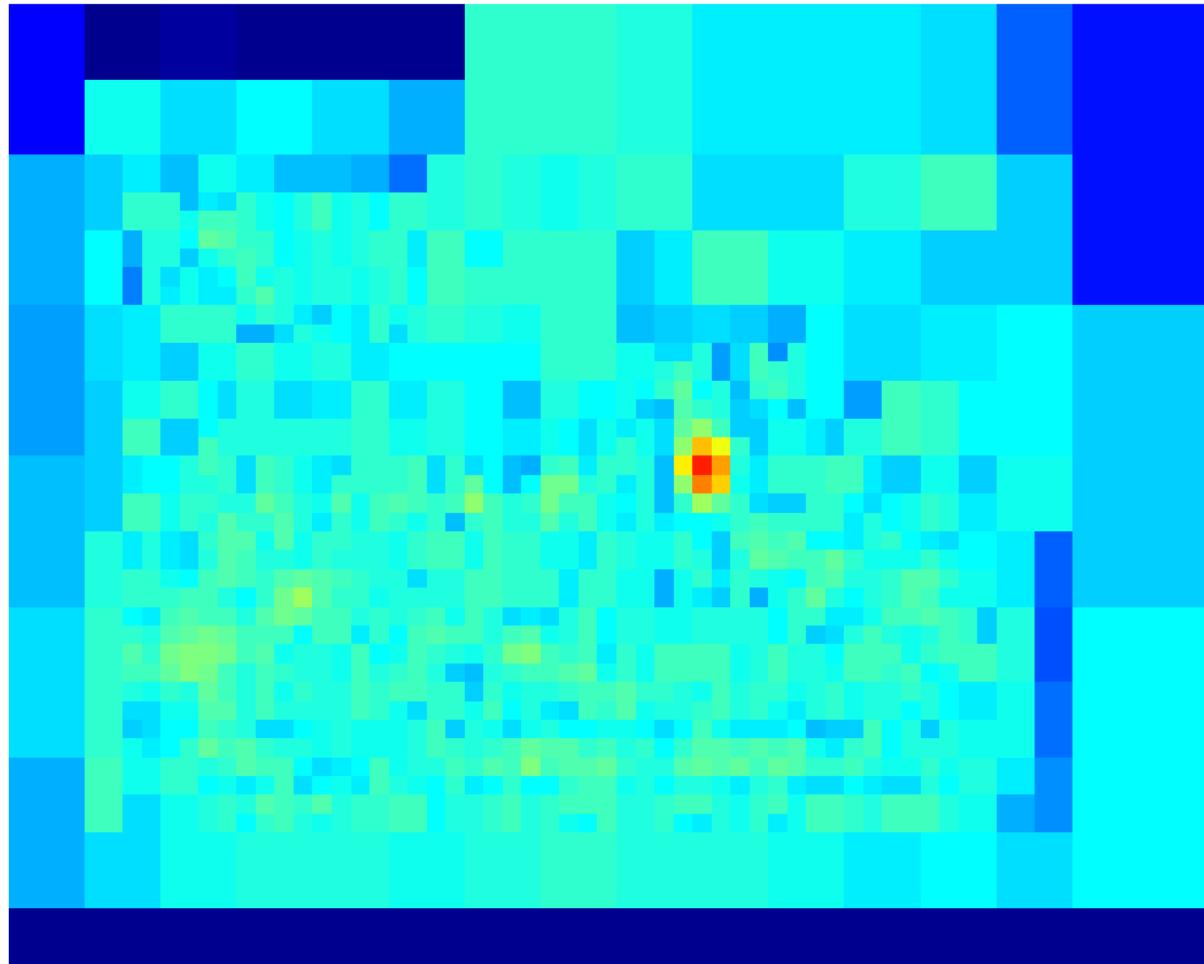
DTBB results, part score bounds @ $p_e = 0.2$



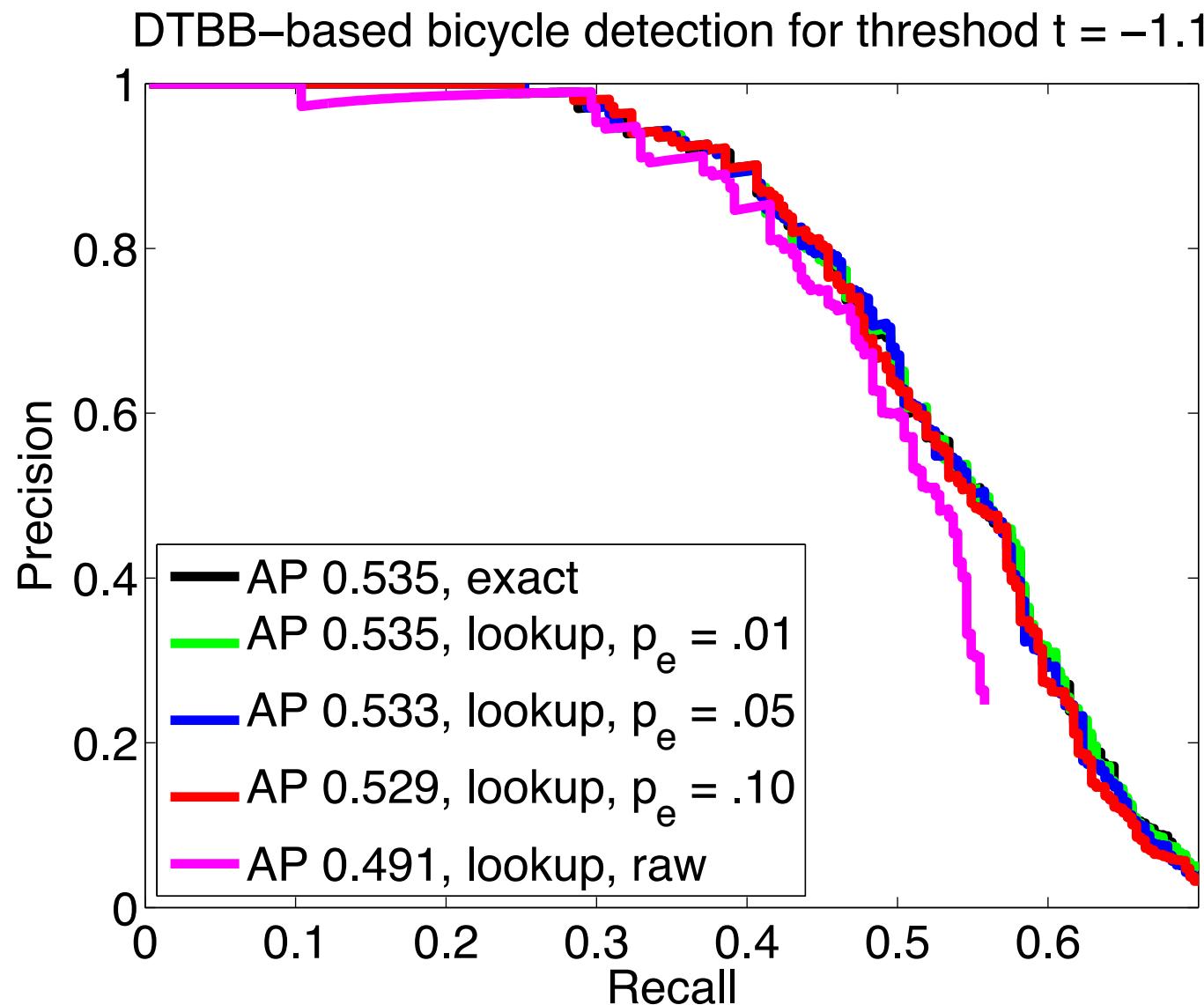
DTBB results, part score bounds @ $p_e = 0.1$



DTBB results, part score bounds @ $p_e = 0.05$



Impact on performance



Speedup results

	GDTs	DTBB	$p_e = 0.05$	$p_e = 0.01$
Part terms	8.35 ± 0.77	1.69 ± 0.18	0.69 ± 0.03	0.69 ± 0.06

Detection with Cascade DPMs (C-DPMs)

$$S_0(x) = 0, \quad \mathcal{I}_0 = [1, N] \times [1, M]$$

$$S_k(x) = S_{k-1}(x) + \max_{x'} (U_p(x') + B_p(x', x))$$

$$\mathcal{I}_k = \{x \in \mathcal{I}_{k-1} : S_{k-1}(x) \geq \theta_k\}$$

Felzenszwalb, Girschick, et al: use PCA-projection of \mathbf{h}, \mathbf{w}

Our work: use quick upper bounds, thresholds for full HOG

	GDTs	C-DPM	$p_e = 0.05$	$p_e = 0.01$
$\theta = -0.5$	8.95 ± 0.82	0.56 ± 0.07	0.19 ± 0.03	0.23 ± 0.04
$\theta = -0.7$	8.95 ± 0.82	0.72 ± 0.09	0.29 ± 0.04	0.36 ± 0.06
$\theta = -1.0$	8.95 ± 0.82	1.04 ± 0.16	0.51 ± 0.10	1.07 ± 0.29

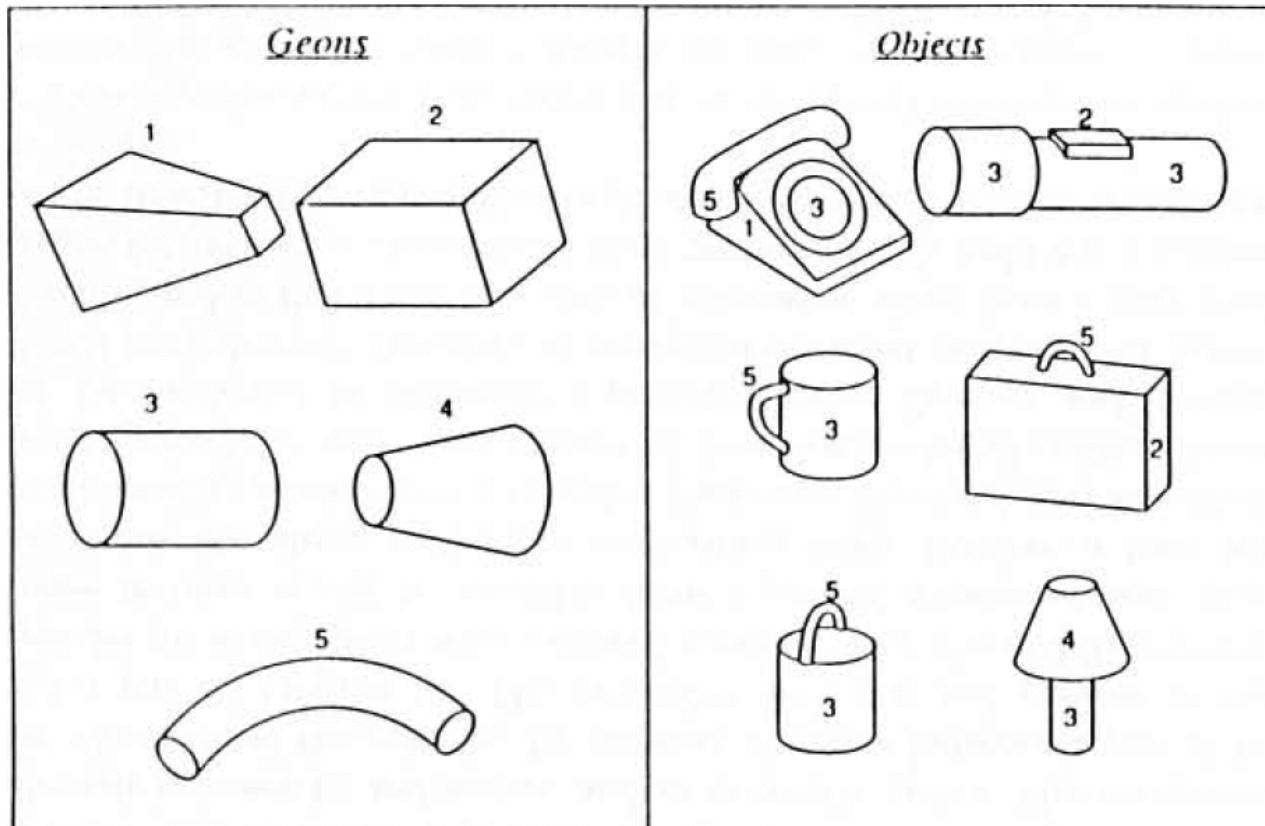


Future work

~10,000 to 30,000



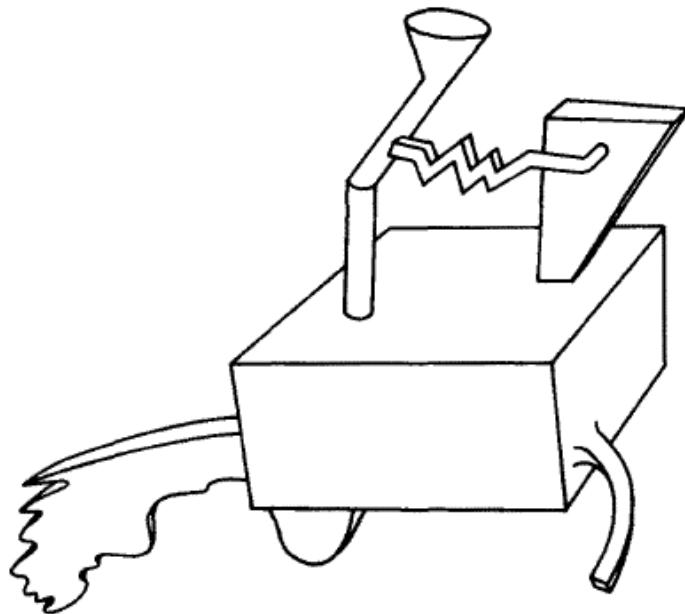
Scalability: objects and their geons



Hypothesis: there is a small number of geometric components that constitute the primitive elements of the object recognition system

Analogy: using letters to form words (compare with Ideograms)

Scalability: Recognition-by-components



- 1) We know that this object is nothing we know**
- 2) We can split this objects into parts that everybody will agree**
- 3) We can see how it resembles something familiar: “a hot dog cart”**

Mid-level vision

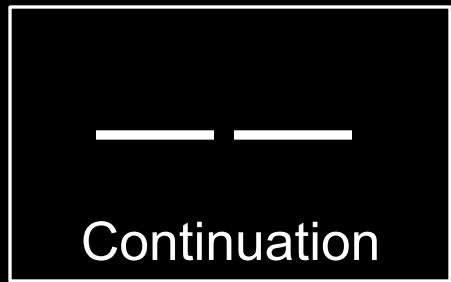
Object-like, but not an object



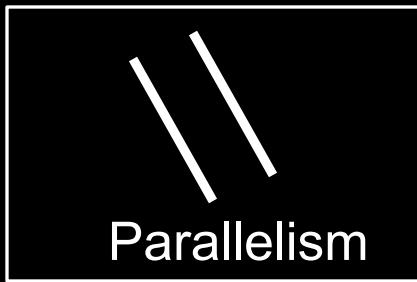
Mid-Level Representations

.....

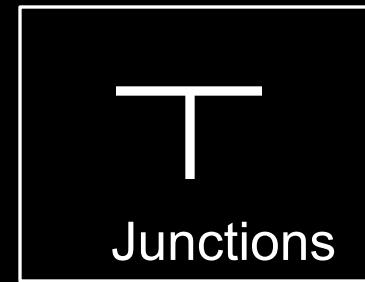
- Mid-level cues



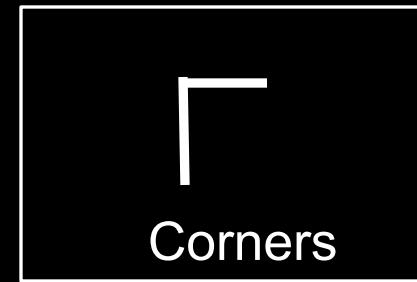
Continuation



Parallelism

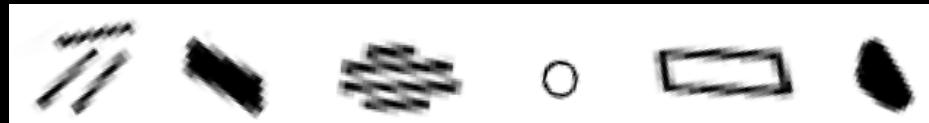


Junctions



Corners

“Tokens” from Vision by D.Marr:



-
- Object parts:

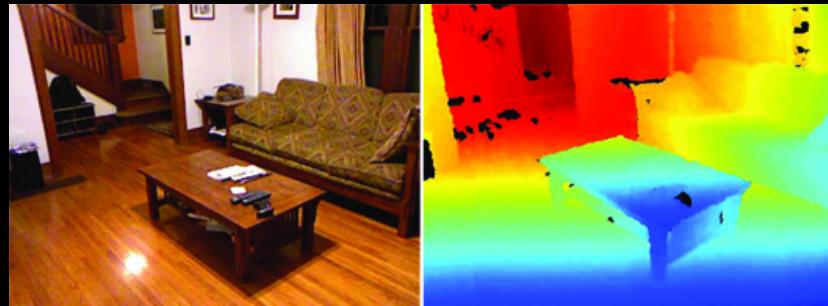


-
- Difficult to hand-engineer → What about learning them?

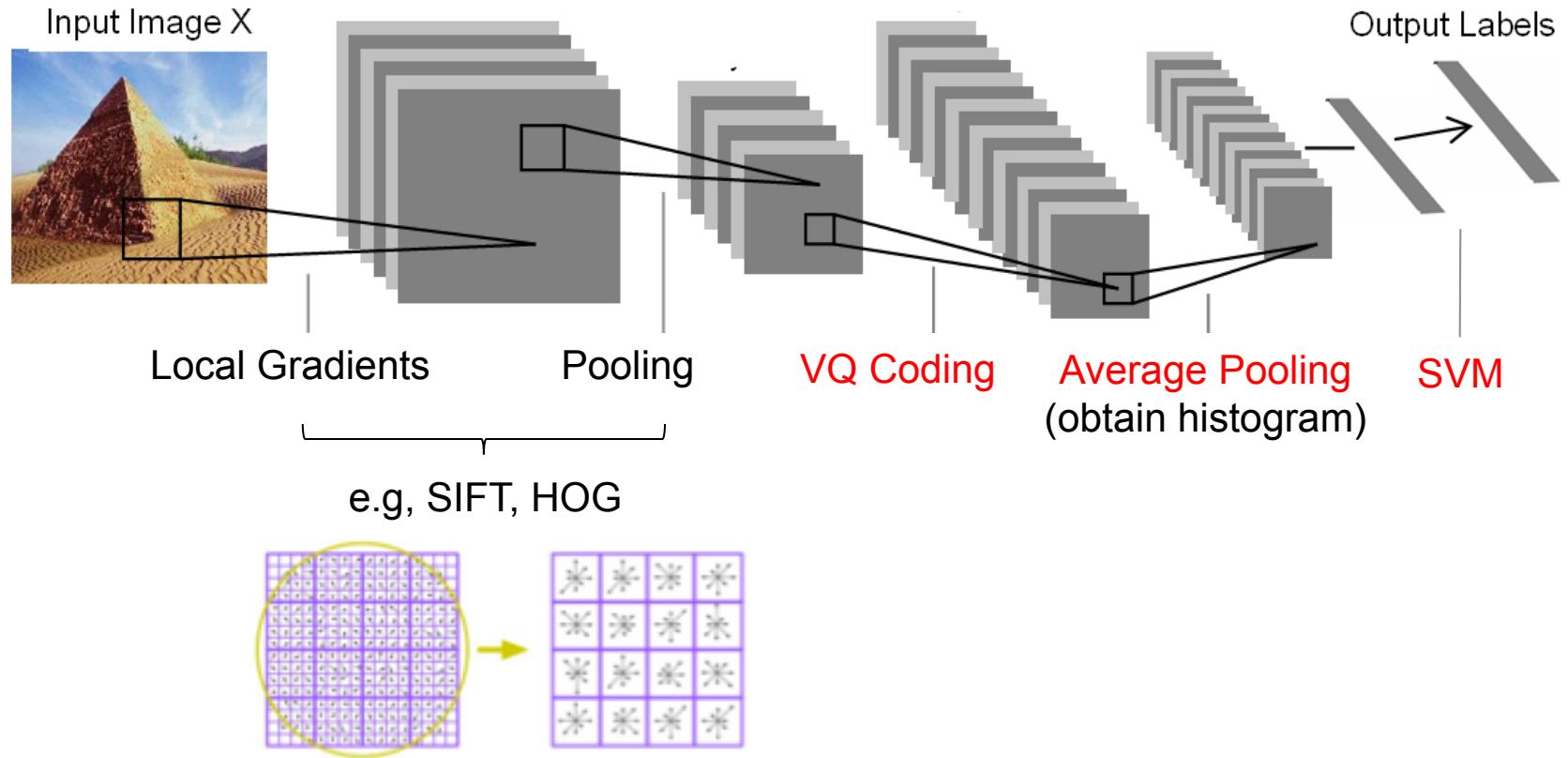
Why Learn Features?

.....

- Better performance
- Other domains (unclear how to hand engineer):
 - Kinect
 - Video
 - Multi spectral
- Feature computation time
 - Dozens of features now regularly used
 - Getting prohibitive for large datasets (10's sec /image)

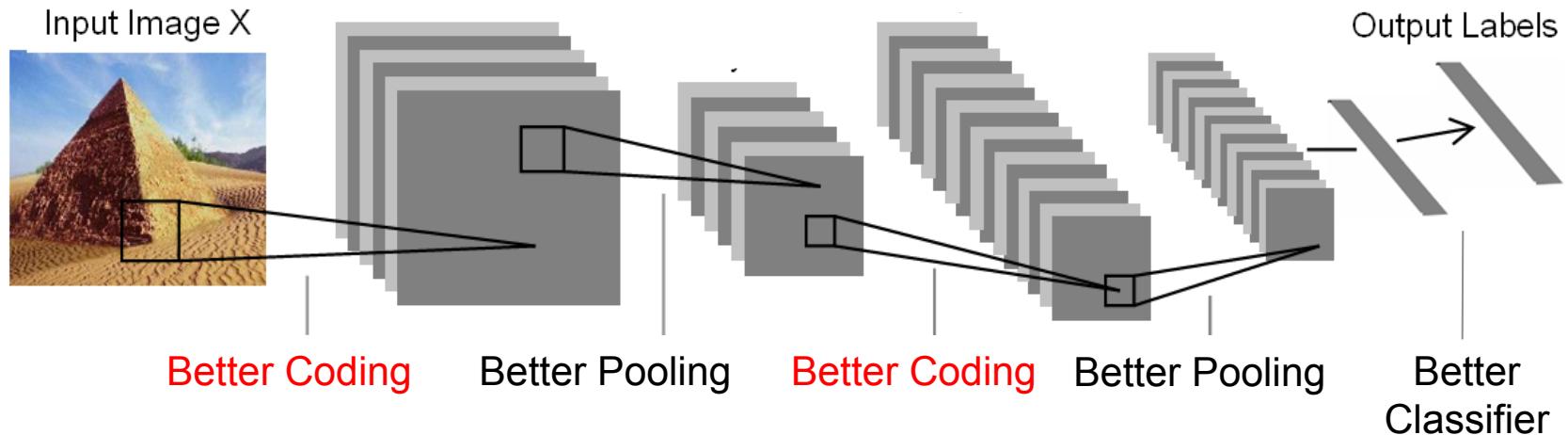


“BoW+SPM” has two coding+pooling layers



SIFT feature itself follows a coding+pooling operation

Develop better coding methods



- Coding: nonlinear mapping data into another feature space
- Better coding methods: **sparse coding**, RBMs, auto-encoders

What is sparse coding

Sparse coding (Olshausen & Field, 1996). Originally developed to explain early visual processing in the brain (edge detection).

Training: given a set of random patches x , learning a dictionary of bases $[\Phi_1, \Phi_2, \dots]$

Coding: for data vector x , solve LASSO to find the sparse coefficient vector a

$$\min_{a, \phi} \sum_{i=1}^m \left\| x_i - \sum_{j=1}^k a_{i,j} \phi_j \right\|^2 + \lambda \sum_{i=1}^m \sum_{j=1}^k |a_{i,j}|$$

Sparse coding: training time

Input: Images x_1, x_2, \dots, x_m (each in \mathbb{R}^d)

Learn: Dictionary of bases $\phi_1, \phi_2, \dots, \phi_k$ (also \mathbb{R}^d).

$$\min_{a, \phi} \sum_{i=1}^m \left\| x_i - \sum_{j=1}^k a_{i,j} \phi_j \right\|^2 + \lambda \sum_{i=1}^m \sum_{j=1}^k |a_{i,j}|$$

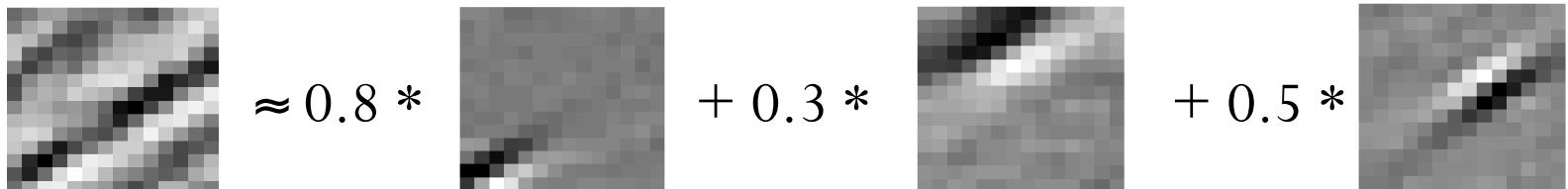
Alternating optimization:

1. Fix dictionary $\phi_1, \phi_2, \dots, \phi_k$, optimize a (a standard LASSO problem)
2. Fix activations a , optimize dictionary $\phi_1, \phi_2, \dots, \phi_k$ (a convex QP problem)

Sparse coding: testing time

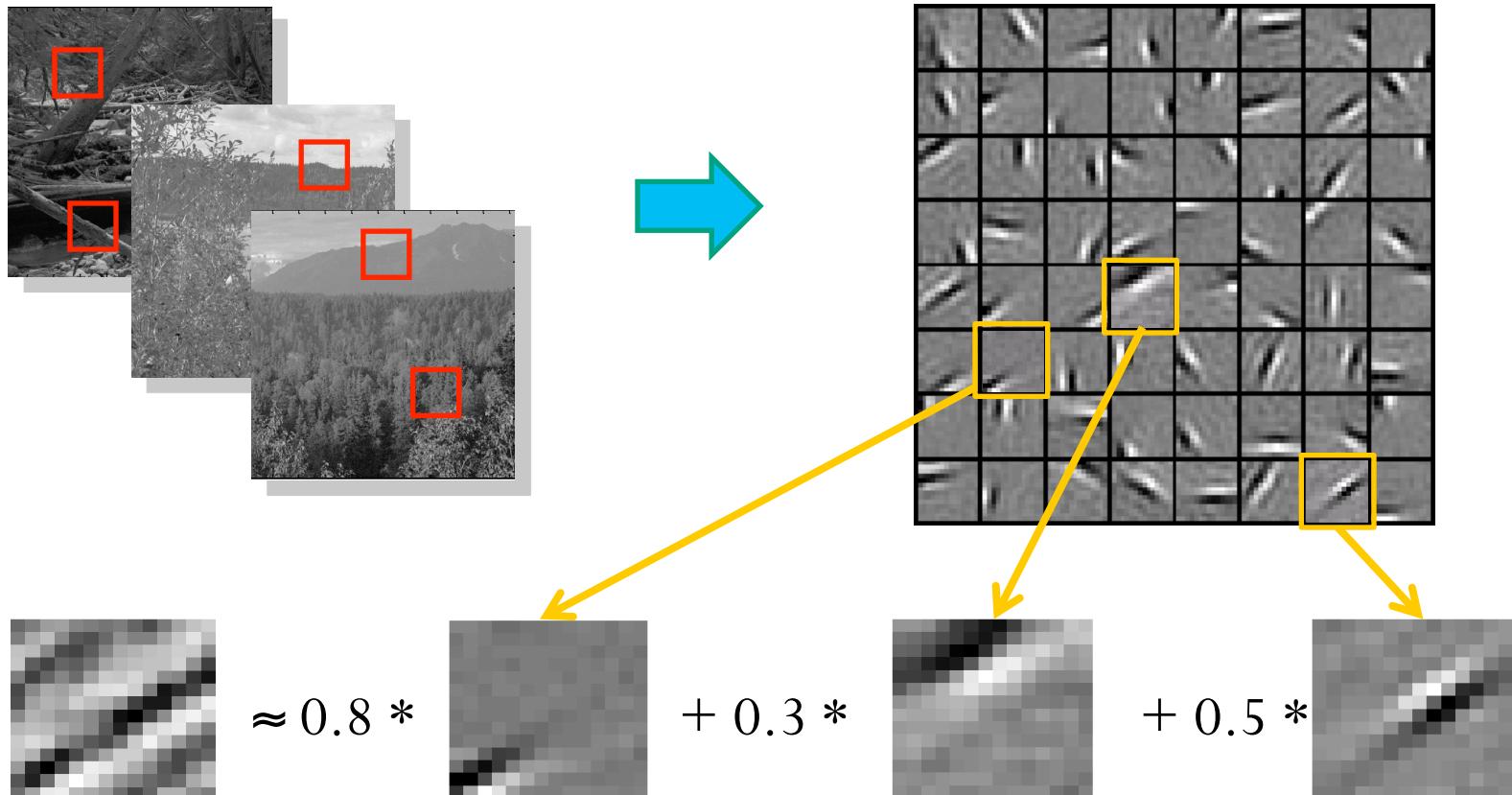
Input: Novel image patch x (in R^d) and previously learned ϕ_i 's
Output: Representation $[a_{i,1}, a_{i,2}, \dots, a_{i,k}]$ of image patch x_i .

$$\min_a \sum_{i=1}^m \left\| x_i - \sum_{j=1}^k a_{i,j} \phi_j \right\|^2 + \lambda \sum_{i=1}^m \sum_{j=1}^k |a_{i,j}|$$



Represent x_i as: $a_i = [0, 0, \dots, 0, \mathbf{0.8}, 0, \dots, 0, \mathbf{0.3}, 0, \dots, 0, \mathbf{0.5}, \dots]$

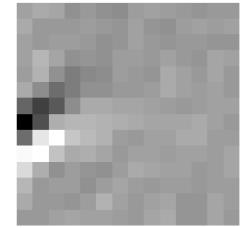
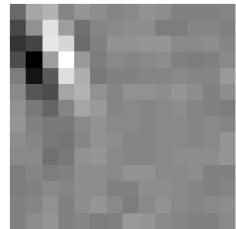
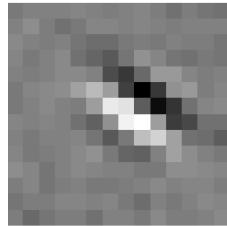
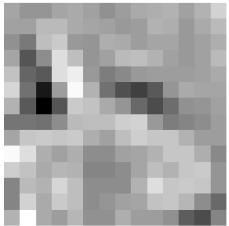
Sparse coding illustration



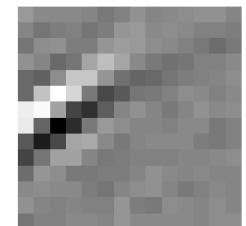
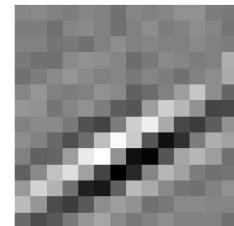
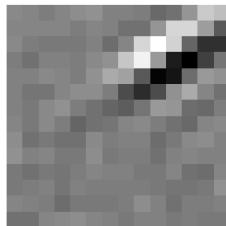
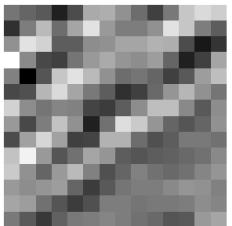
$$[a_1, \dots, a_{64}] = [0, 0, \dots, 0, \mathbf{0.8}, 0, \dots, 0, \mathbf{0.3}, 0, \dots, 0, \mathbf{0.5}, 0]$$

(feature representation)

More examples

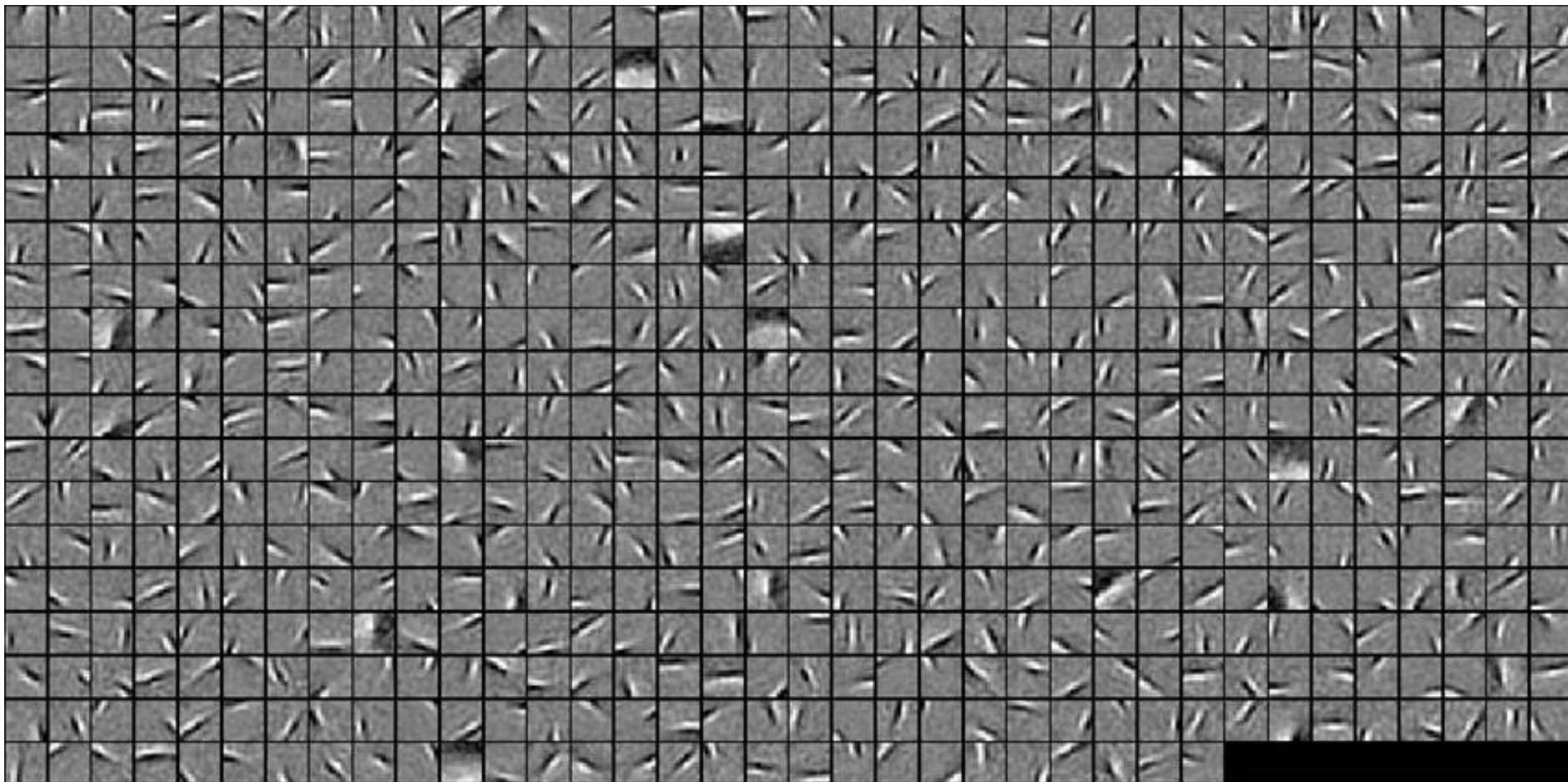


Represent as: [0, 0, ..., 0, 0.6, 0, ..., 0, 0.8, 0, ..., 0, 0.4, ...]



Represent as: [0, 0, ..., 0, 1.3, 0, ..., 0, 0.9, 0, ..., 0, 0.3, ...]

Sparse coding basis



Recap of sparse coding for feature learning

Training time

Input: Images $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (each in $\mathbb{R}^{n \times n}$)

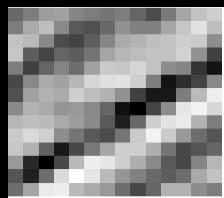
Learn: Dictionary of bases $\phi_1, \phi_2, \dots, \phi_k$ (also $\mathbb{R}^{n \times n}$).

$$\min_{a, \phi} \sum_{i=1}^m \left(\left\| x^{(i)} - \sum_{j=1}^k a_j^{(i)} \phi_j \right\|^2 - \lambda \sum_{j=1}^k |a_j^{(i)}| \right)$$

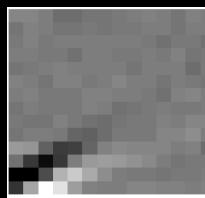
Input: Novel image x (in $\mathbb{R}^{n \times n}$) and previously learned ϕ_i 's.
Output: Representation $[a_1, a_2, \dots, a_k]$ of image x .

Test time

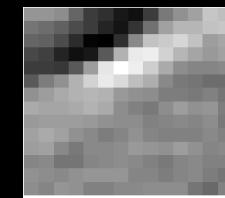
$$\min_a \left\| x - \sum_{j=1}^k a_j \phi_j \right\|^2 + \lambda \sum_{j=1}^k |a_j|$$



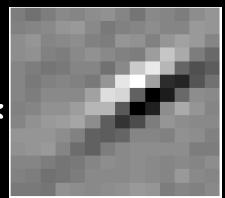
$$\approx 0.8 *$$



$$+ 0.3 *$$



$$+ 0.5 *$$



$$x$$

$$\approx 0.8 *$$

$$\phi_{36}$$

$$+ 0.3 *$$

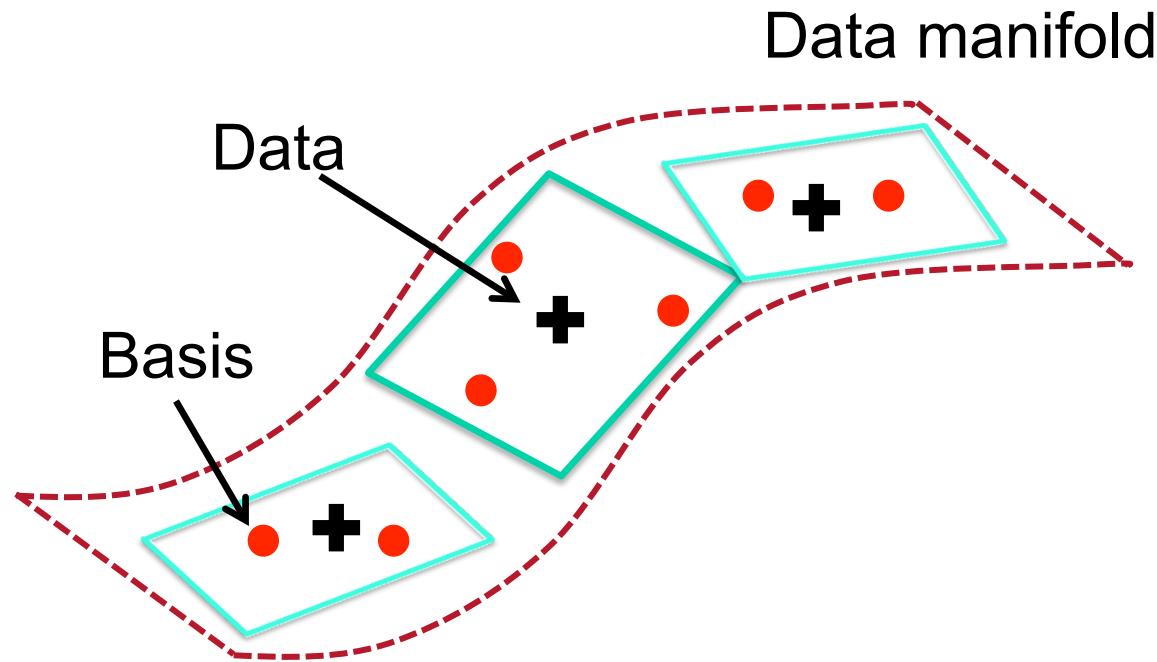
$$\phi_{42}$$

$$+ 0.5 *$$

$$\phi_{63}$$

Represent as: [0, 0, ..., 0, **0.8**, 0, ..., 0, **0.3**, 0, ..., 0, **0.5**, ...]

A geometric view to sparse coding

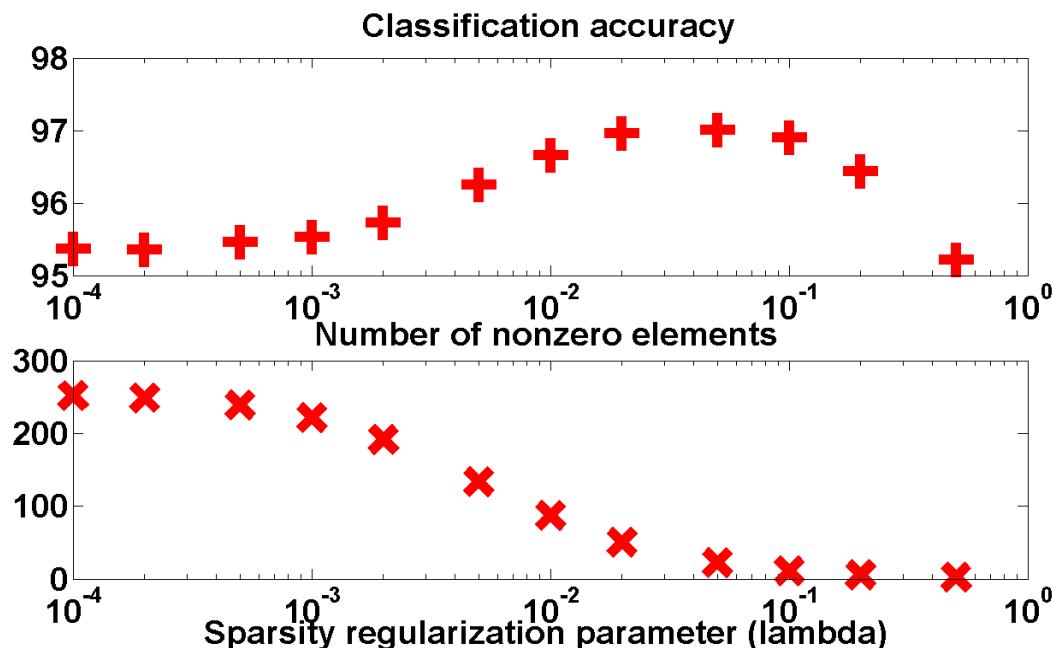


- Each basis is somewhat like a pseudo data point – “**anchor point**”
- **Sparsity**: each datum is a sparse combination of neighbor anchors.
- The coding scheme explores the **manifold structure** of data.

MNIST Experiment: Classification using SC

$$\min_{a, \phi} \sum_{i=1}^m \left\| x_i - \sum_{j=1}^k a_{i,j} \phi_j \right\|^2 + \lambda \sum_{i=1}^m \sum_{j=1}^k |a_{i,j}|$$

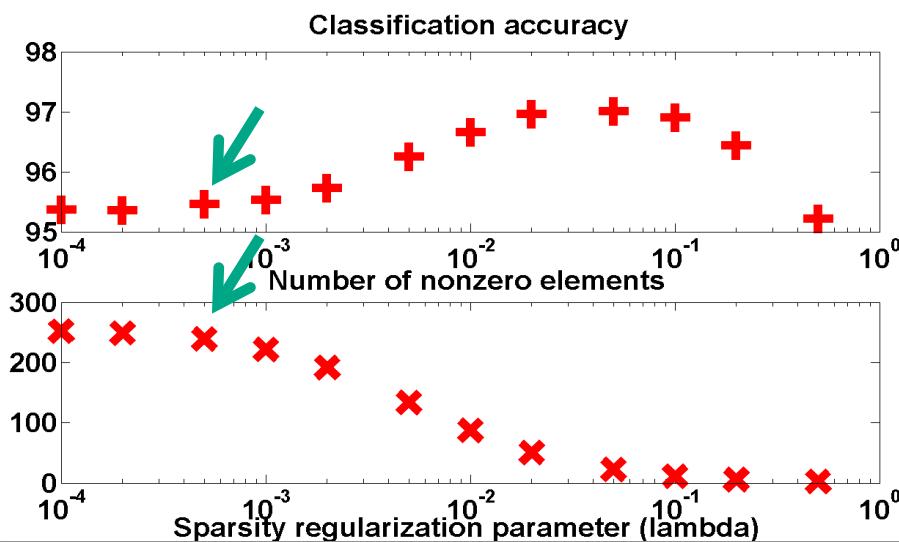
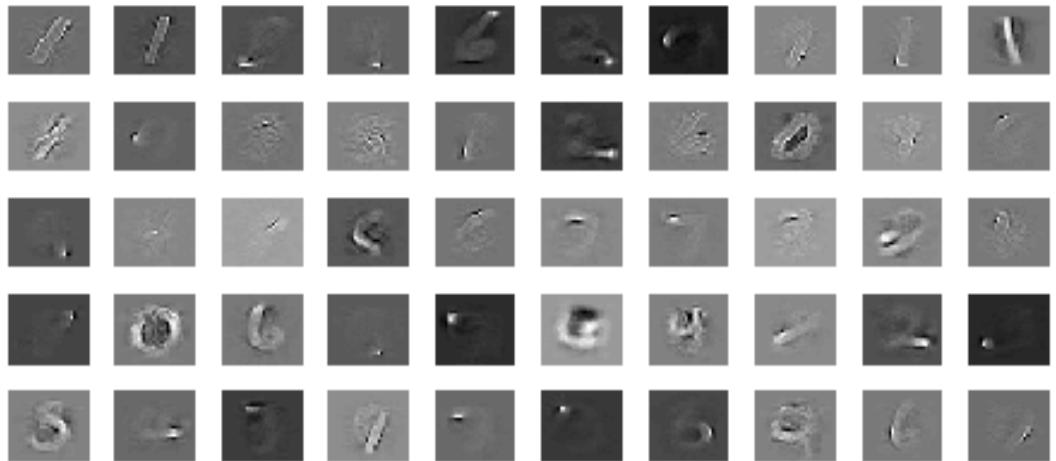
Try different values



- 60K training, 10K for test
- Let k=512
- Linear SVM on sparse codes

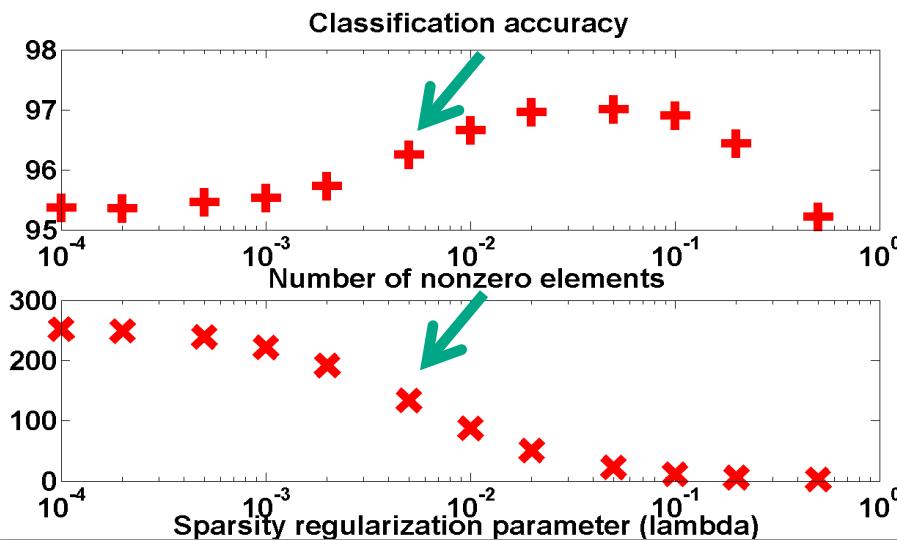
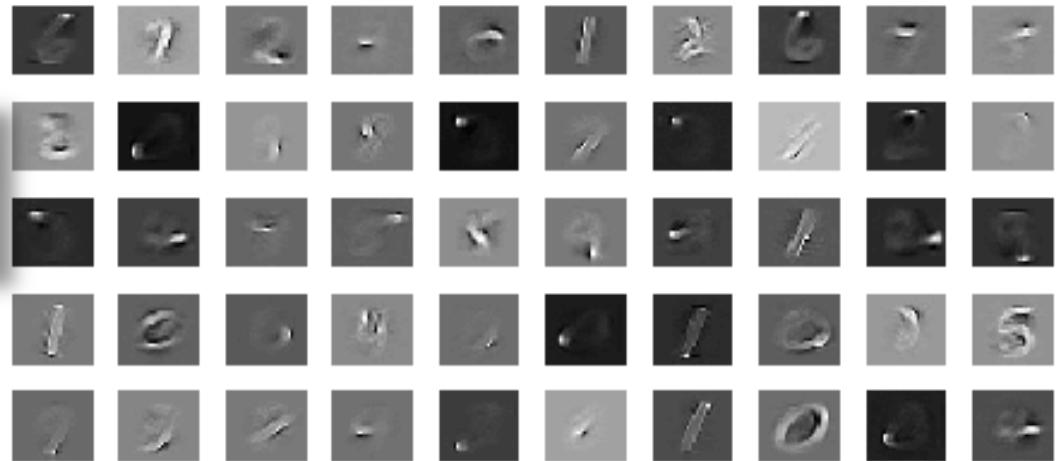
MNIST Experiment: Lambda = 0.0005

Each basis is like a part or direction.



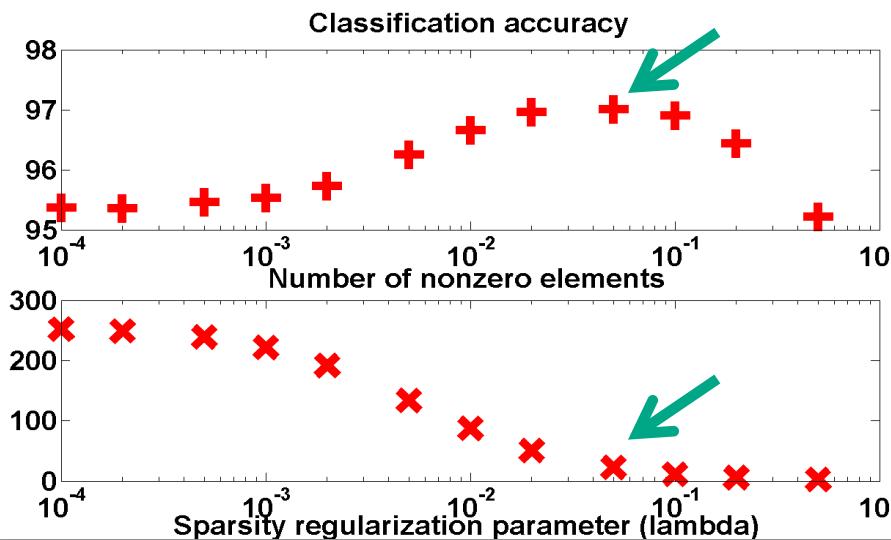
MNIST Experiment: Lambda = 0.005

Again, each basis is like
a part or direction.



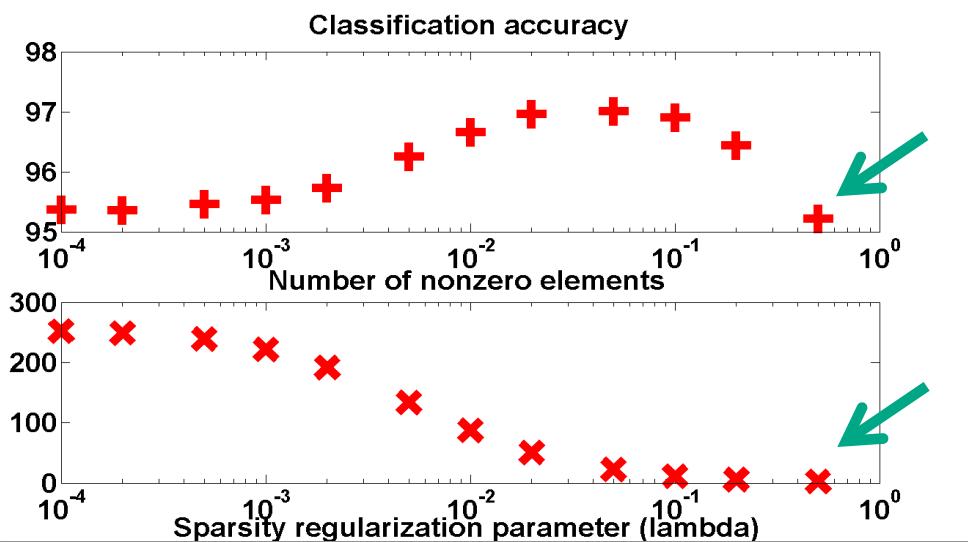
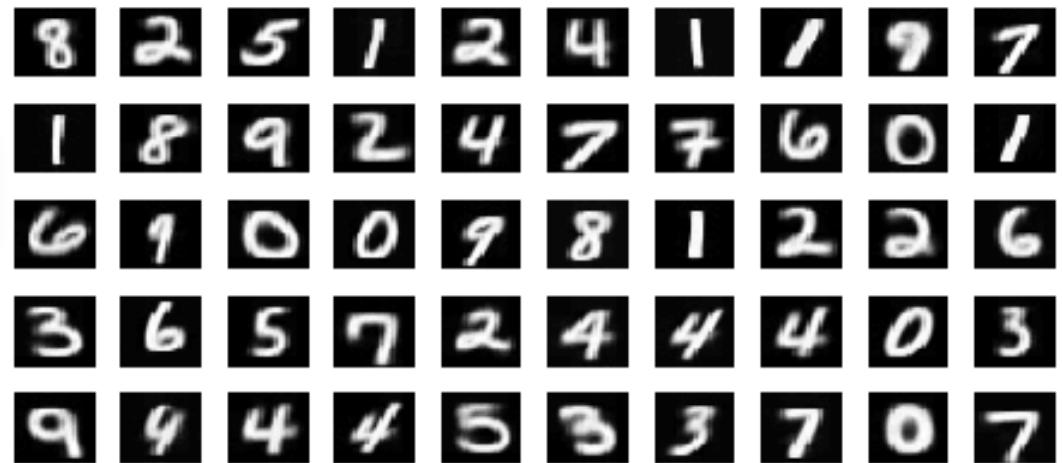
MNIST Experiment: Lambda = 0.05

Now, each basis is more like a digit !

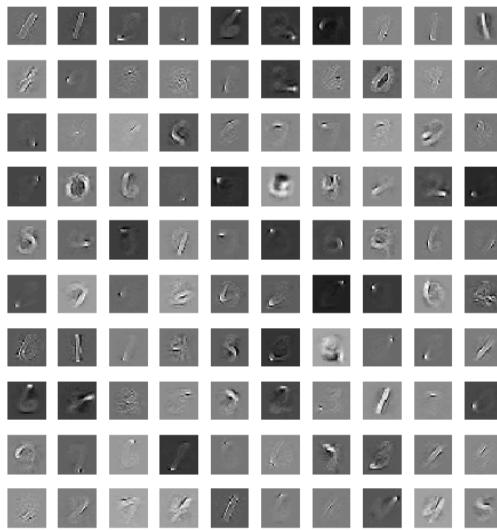


MNIST Experiment: Lambda = 0.5

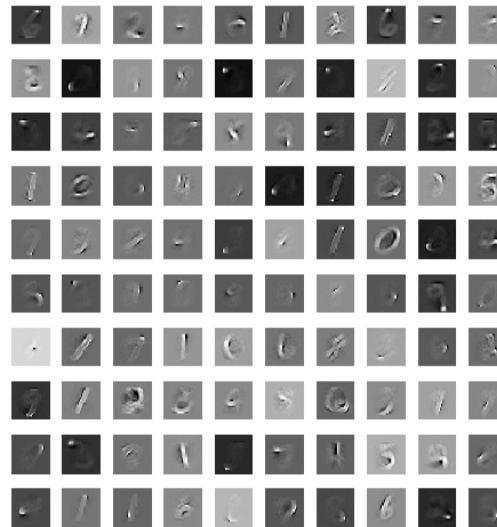
Like clustering now!



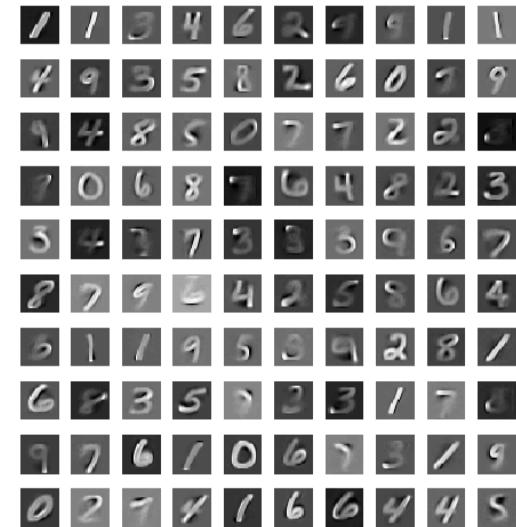
Geometric view of sparse coding



Error: 4.54%



Error: 3.75%



Error: 2.64%

- When SC achieves the best classification accuracy, the learned bases are like digits – **each basis has a clear local class association.**
- Implication: exploring data geometry may be useful for classification.

K-means vs. sparse coding

Intuition: “Soft” version of k-means (membership in multiple clusters).

K-means

Centroid 1

Centroid 2

Centroid 3

Represent as:

$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Sparse

Basis ϕ_1

Basis ϕ_2

Basis ϕ_3

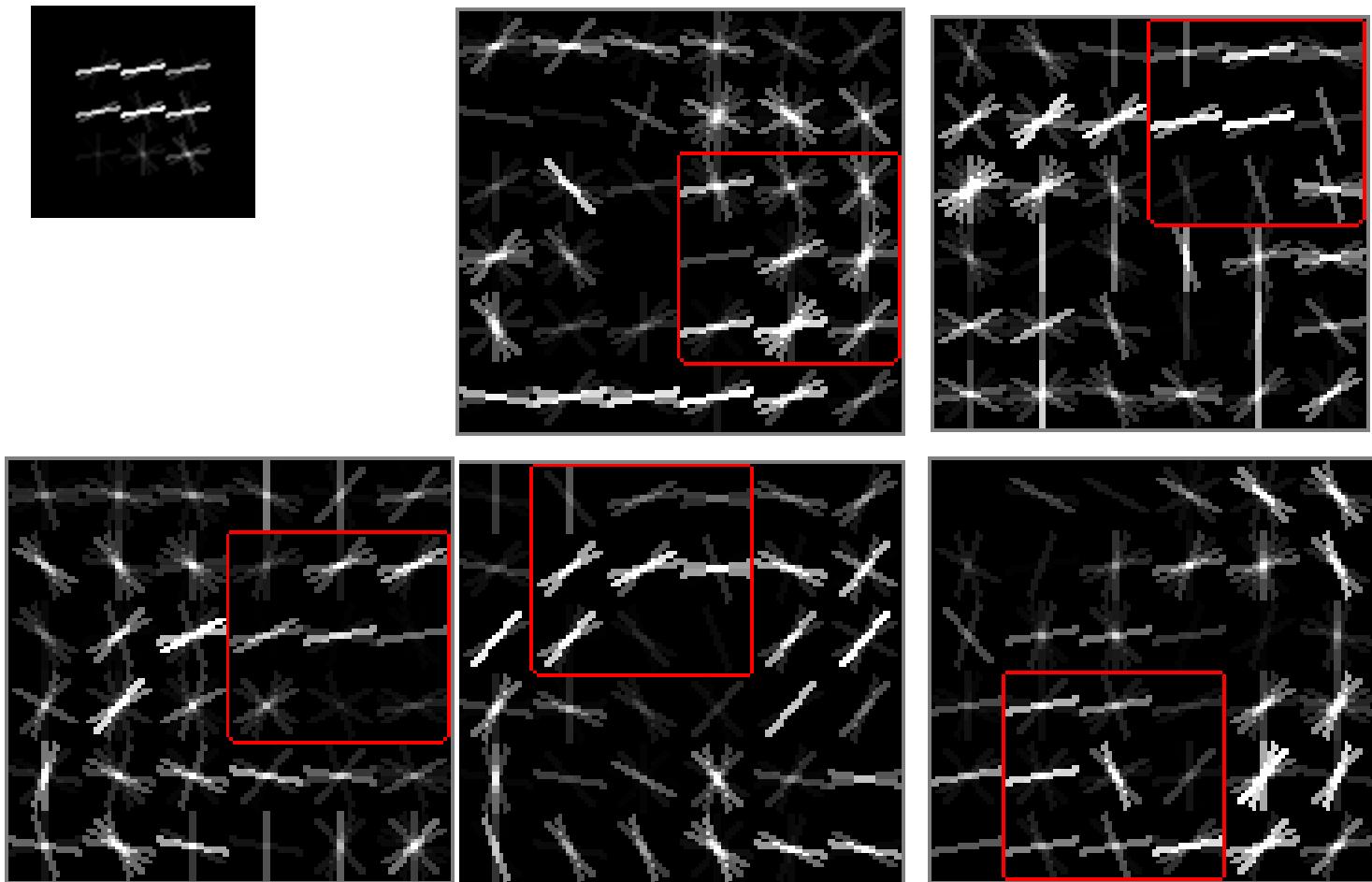
$$\approx 0.3\phi_1 + 0.1\phi_2 + 0.6\phi_3$$

Represent as:

$$\begin{bmatrix} 0.3 \\ 0.1 \\ 0.6 \end{bmatrix}$$

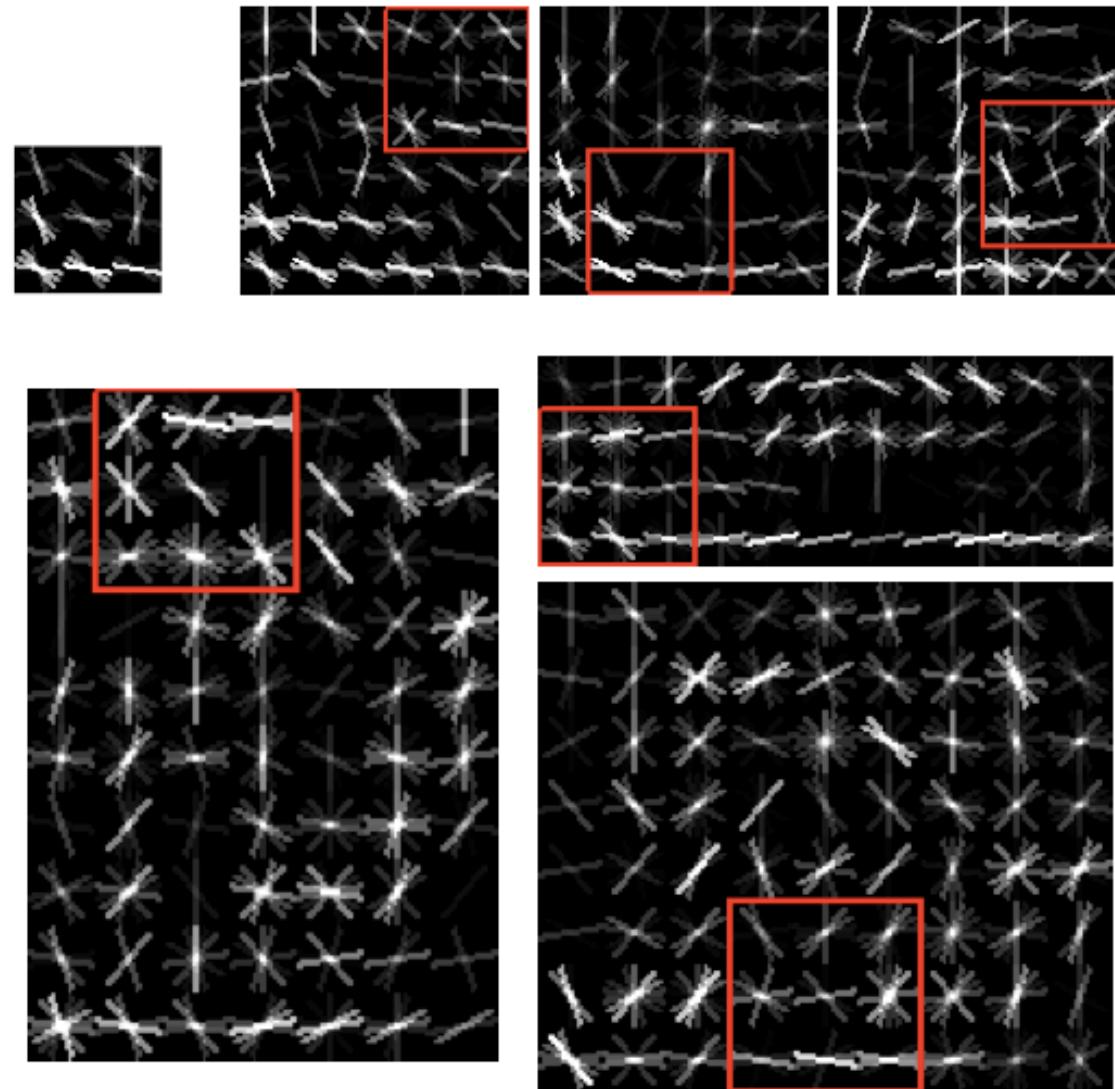
Dealing with translations

Motivation



Dealing with translations

Motivation



Sparse Coding

$$(P1) \quad \min_{\mathbf{b}, \alpha_1 \dots \alpha_P} \sum_{p=1}^P \|\mathbf{w}_p - \mathbf{b}\alpha_p\|_2$$

s.t. $\|\alpha_p\|_0 \leq L, \quad p = 1, \dots, P$ Notice the $\| \cdot \|_0$ norm
 $\|\mathbf{b}_b\|_2 \leq 1, \quad b = 1, \dots, B$

Olshausen and Field, Sparse coding with an overcomplete basis, Vision Research 1997
Mallat and Zhang, 'Matching Pursuits with time-frequency dictionaries, IEEE Trans. Sig. Proc, 1993

Shift Invariant Sparse Coding

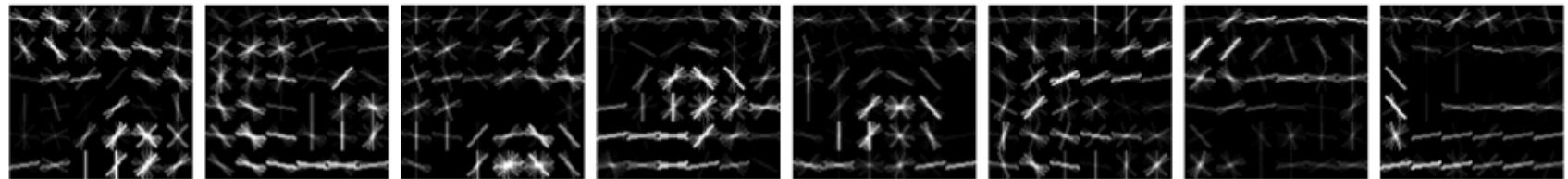
$$(P2) \quad \min_{\mathbf{k}, \alpha_1 \dots \alpha_P} \sum_{p=1}^P \|\mathbf{w}_p - \mathbf{b}\alpha_p\|_2$$

s.t. $\|\alpha_p\|_0 \leq L, \quad p = 1, \dots, P$

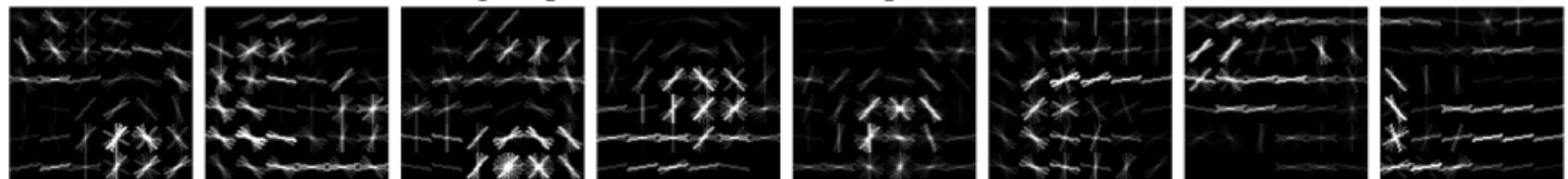
$\|\mathbf{k}_i\|_2 \leq 1, \quad i = 1, \dots, K$

$\mathbf{b}_b = \mathbf{k}_b^{v_b, h_b}, \quad b = 1, \dots, B$

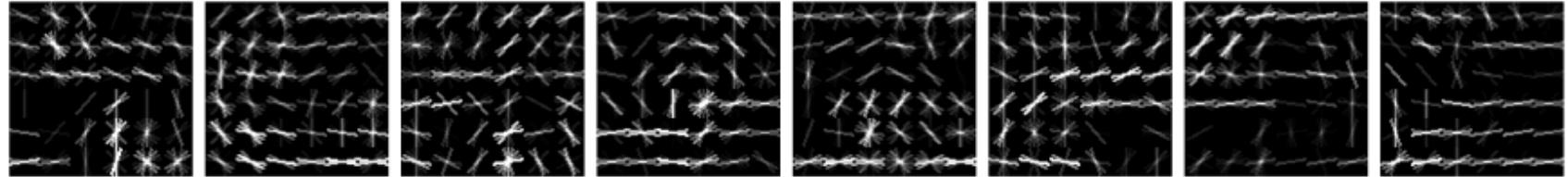
Shufflet-based reconstructions (SISC vs SC)



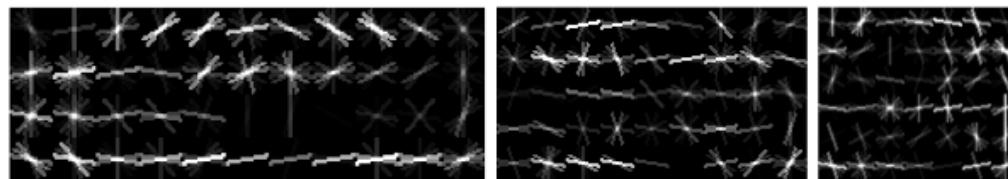
Original part filters for the first component of class 'car'



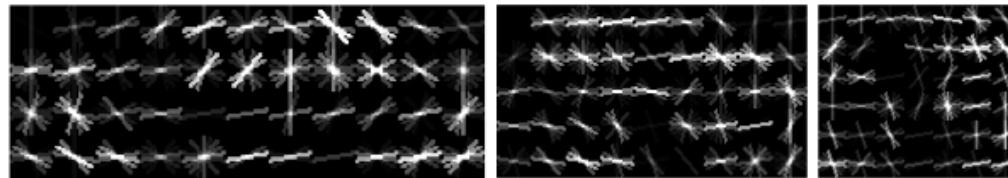
SISC-based reconstructions



Sparse Coding-based reconstructions



Original root filters for components 1, 3, 5 of class 'car'



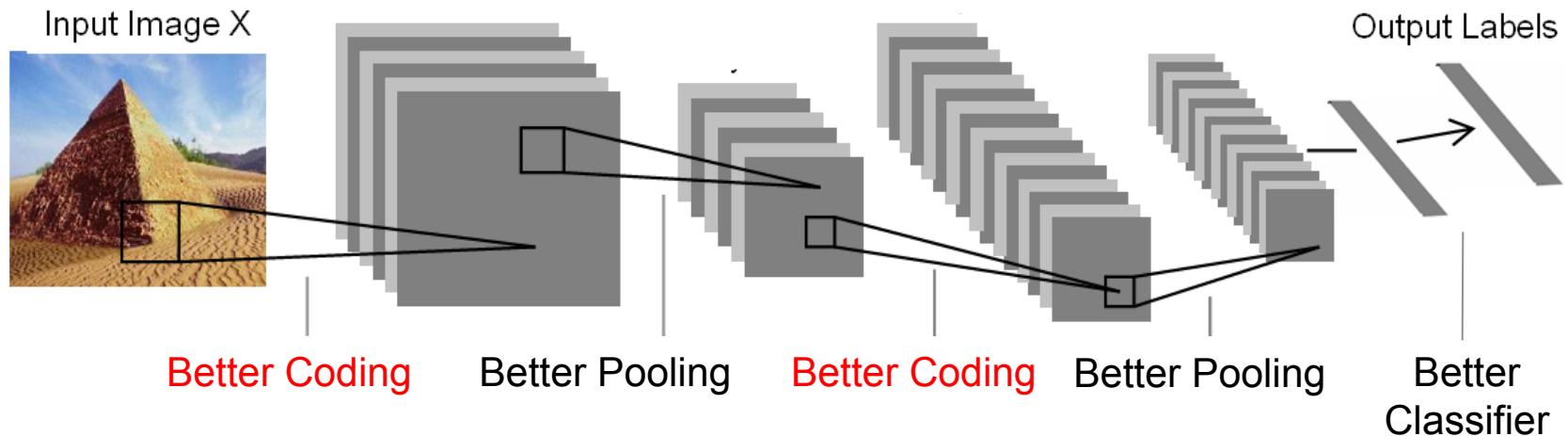
SISC-based approximations

Timing results

	Baseline	DTBB	Lookup- $p_e = 0.5$	Lookup- $p_e = 0.1$	Shufflets- $p_e = 0.5$	Shufflets $p_e = 0.1$
Part terms	2.26 ± 0.77	1.69 ± 0.18	0.69 ± 0.03	0.69 ± 0.06	0.28 ± 0.05	0.28 ± 0.05
$\theta = -0.5$	0.41 ± 0.06	0.21 ± 0.06	0.47 ± 0.11	1.04 ± 0.25	0.54 ± 0.22	0.64 ± 0.22
Sum	2.67 ± 0.83	1.90 ± 0.23	1.17 ± 0.12	1.74 ± 0.32	0.83 ± 0.27	0.93 ± 0.27
$\theta = -0.7$	0.41 ± 0.06	0.42 ± 0.10	1.00 ± 0.23	2.33 ± 0.65	0.74 ± 0.25	0.83 ± 0.26
Sum	2.67 ± 0.83	2.10 ± 0.24	1.70 ± 0.27	3.00 ± 0.71	1.02 ± 0.30	1.13 ± 0.31
$\theta = -1.0$	0.41 ± 0.06	1.31 ± 0.31	3.80 ± 0.90	9.40 ± 2.70	0.87 ± 0.24	1.01 ± 0.27
Sum	2.67 ± 0.83	3.00 ± 0.42	4.50 ± 1.02	10.01 ± 2.82	1.16 ± 0.29	1.29 ± 0.52

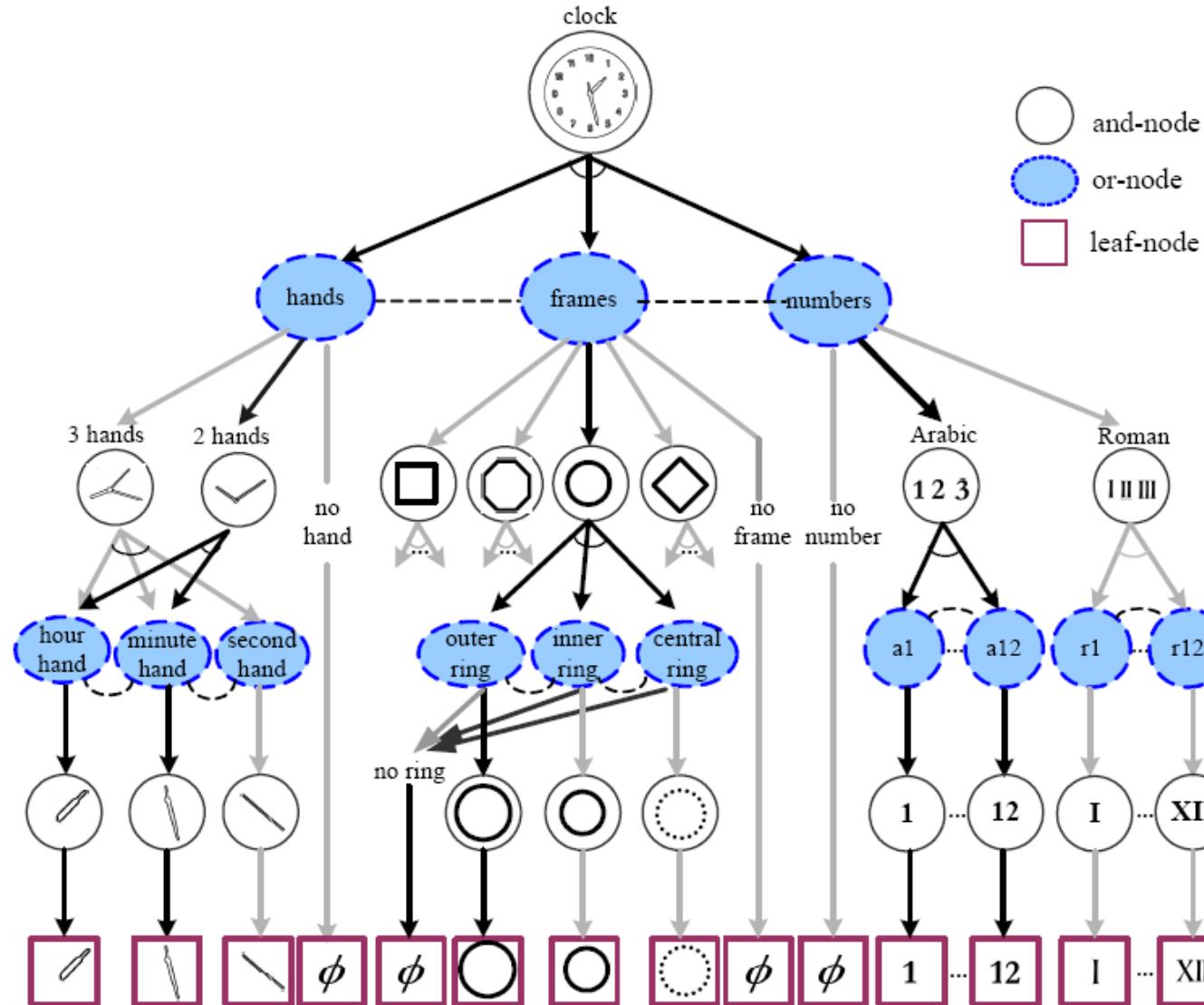
	Baseline	Cascade-DPM	Shufflet cascade, $p_e = 0.1$	Shufflet cascade, $p_e = 0.5$
$\theta = -0.5$	2.67 ± 0.82	0.56 ± 0.07	0.22 ± 0.03	0.17 ± 0.04
$\theta = -0.7$	2.67 ± 0.82	0.72 ± 0.09	0.23 ± 0.04	0.26 ± 0.06
$\theta = -1.0$	2.67 ± 0.82	1.04 ± 0.16	0.46 ± 0.10	0.57 ± 0.29

Fast detection with hierarchical models



Stochastic Grammar of Images

Zhu and Mumford, 'Quest for a Generic Grammar of Images', 2007



Hierarchies in vision and speech

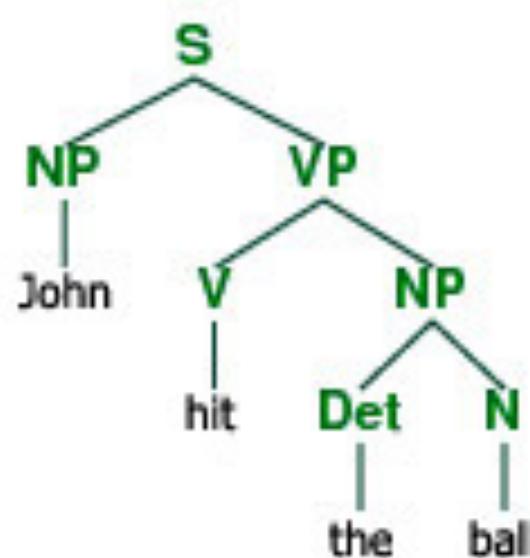
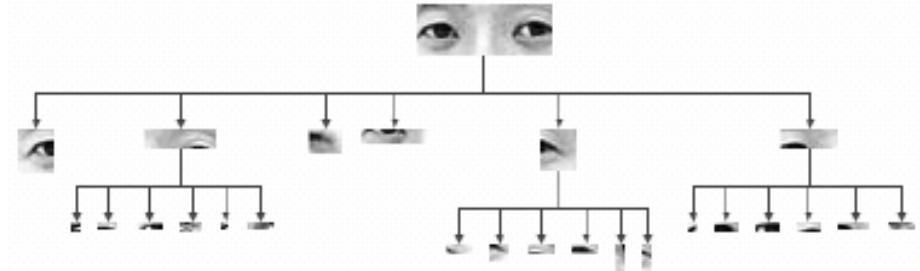
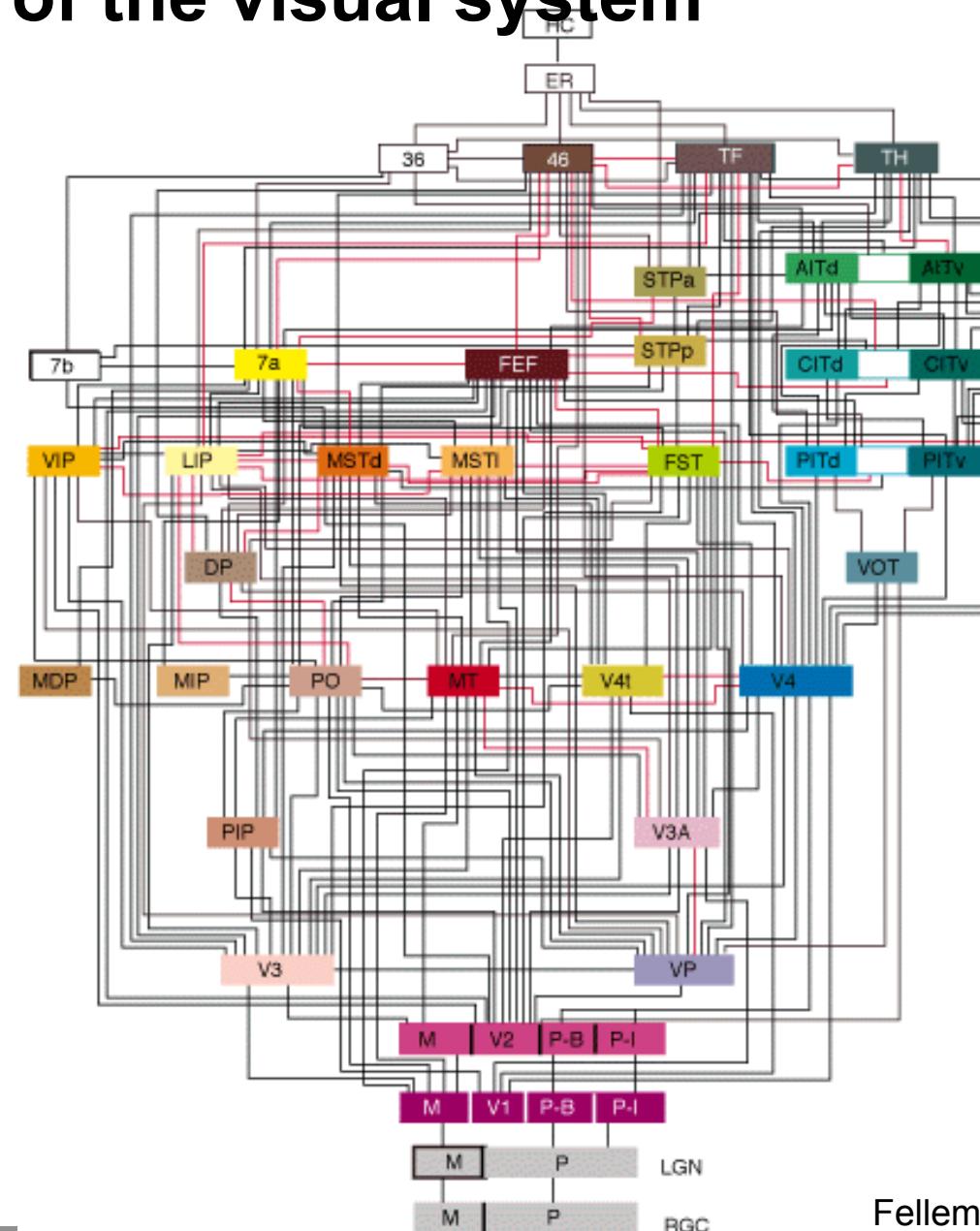
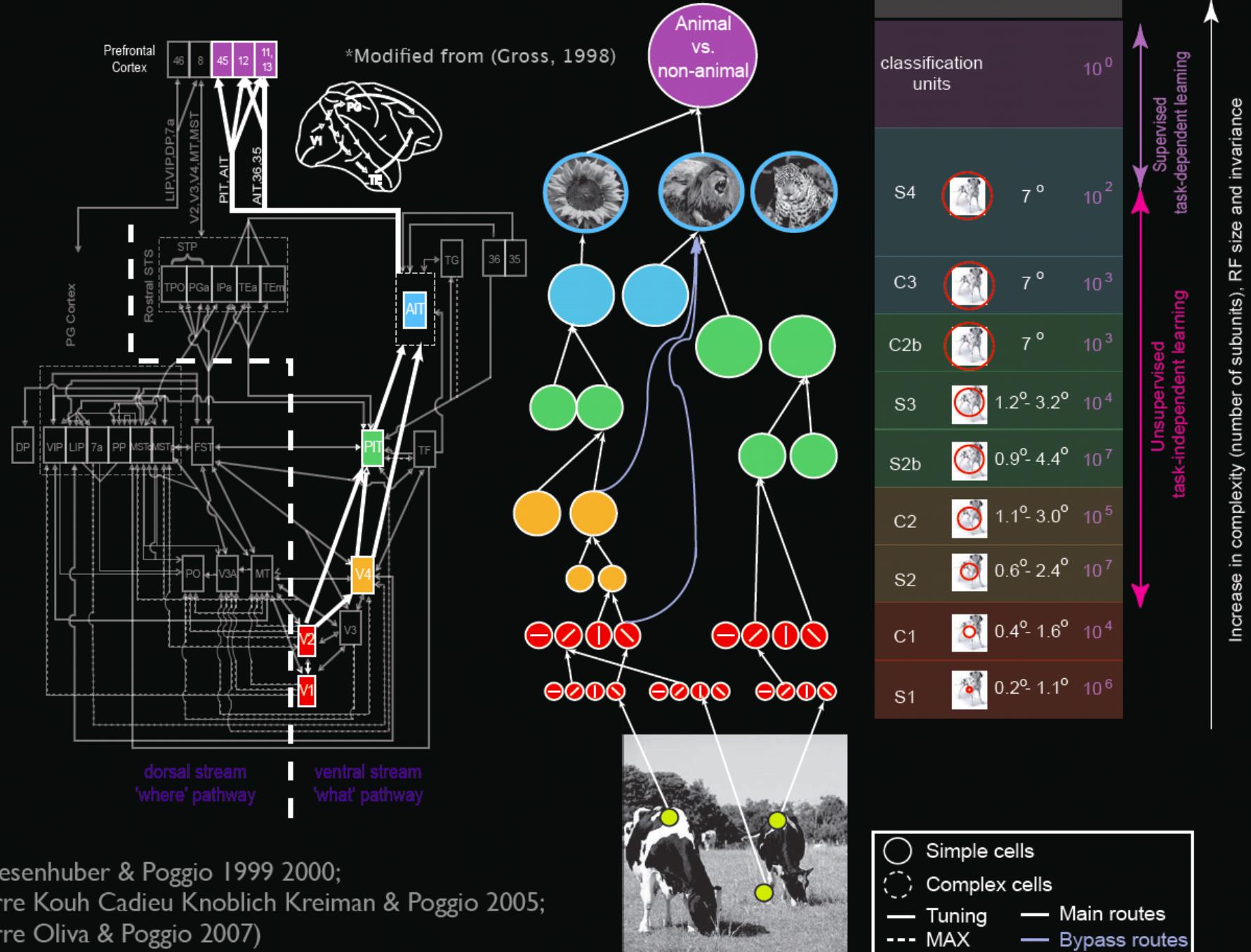
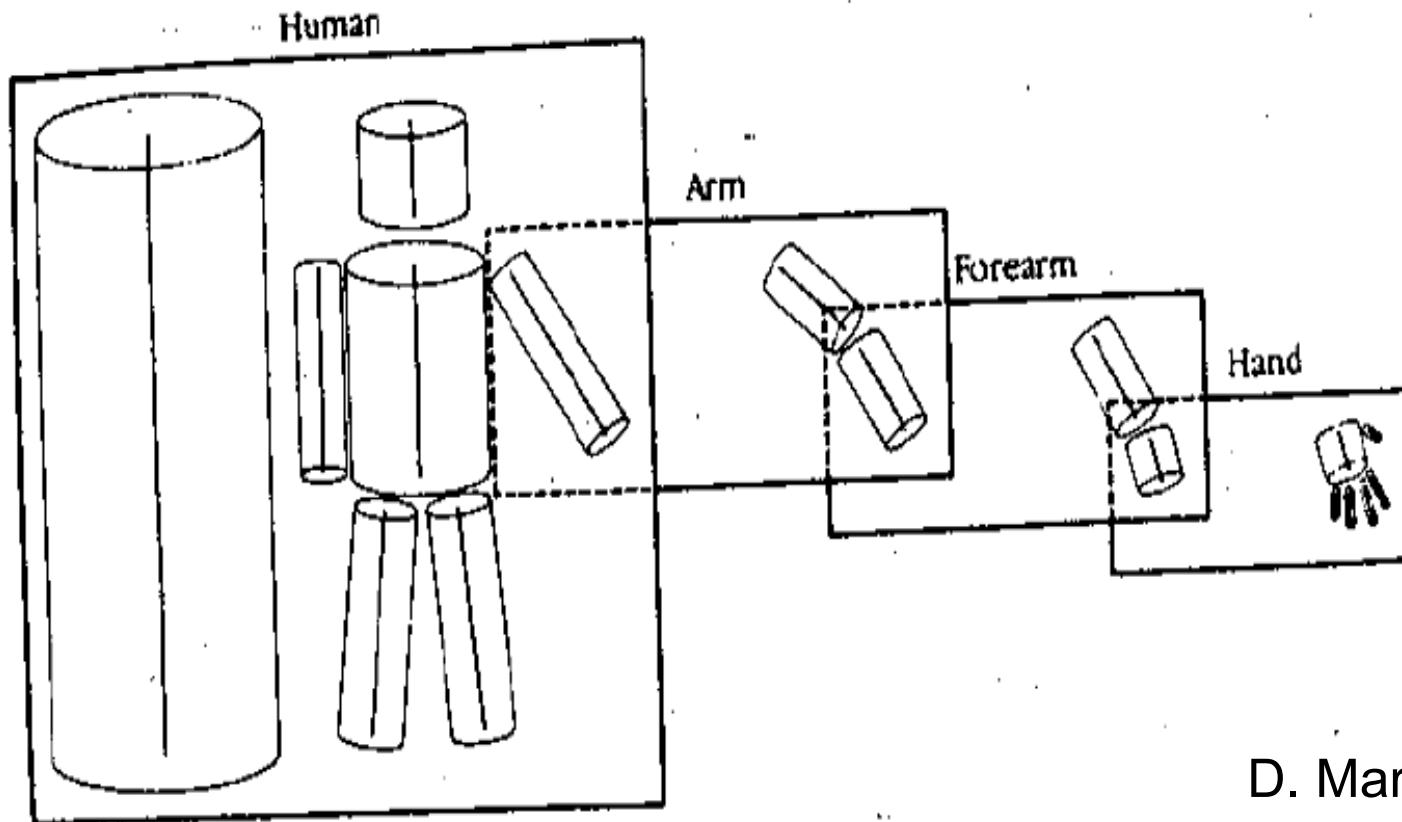


Diagram of the visual system



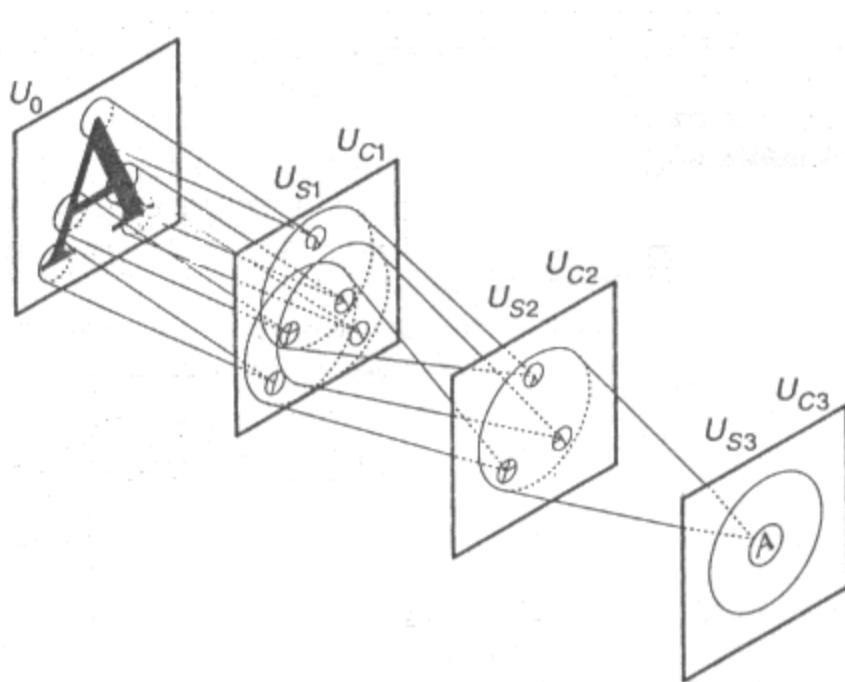


D. Marr's cylinders, 1978



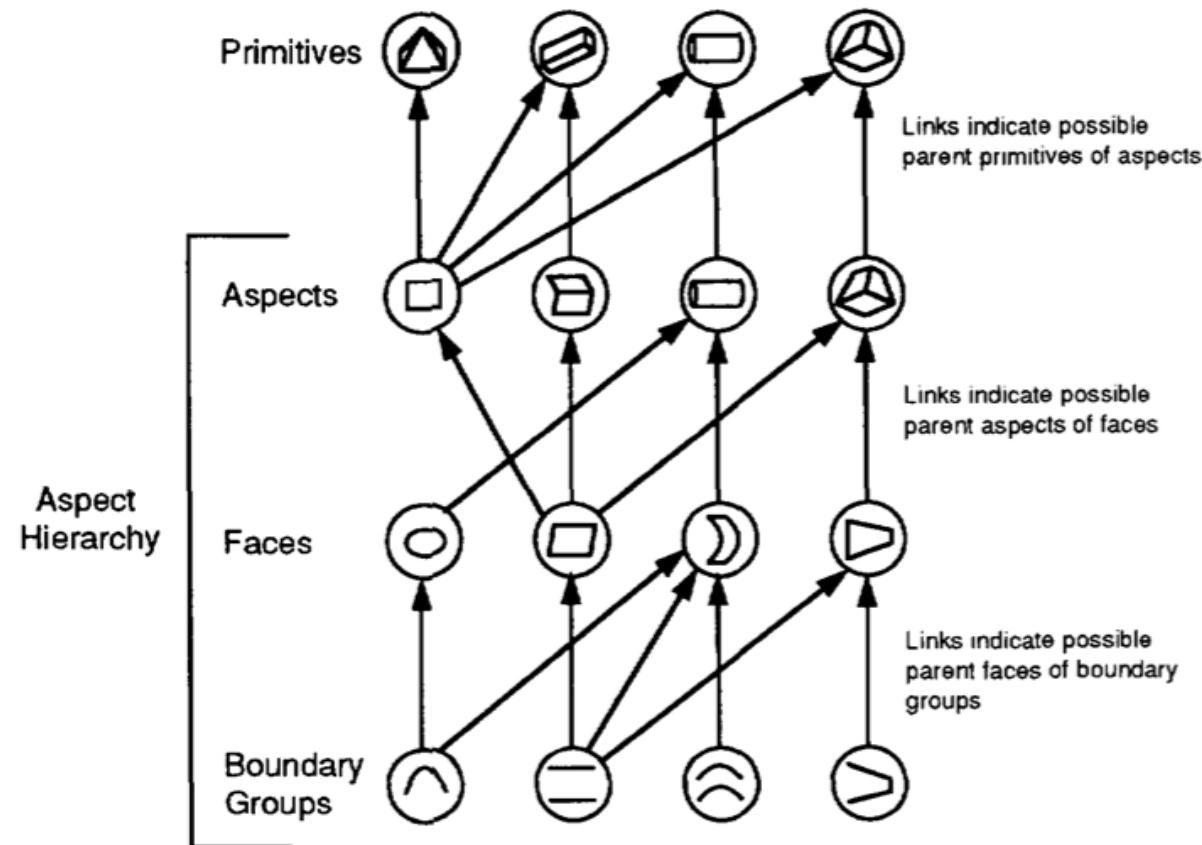
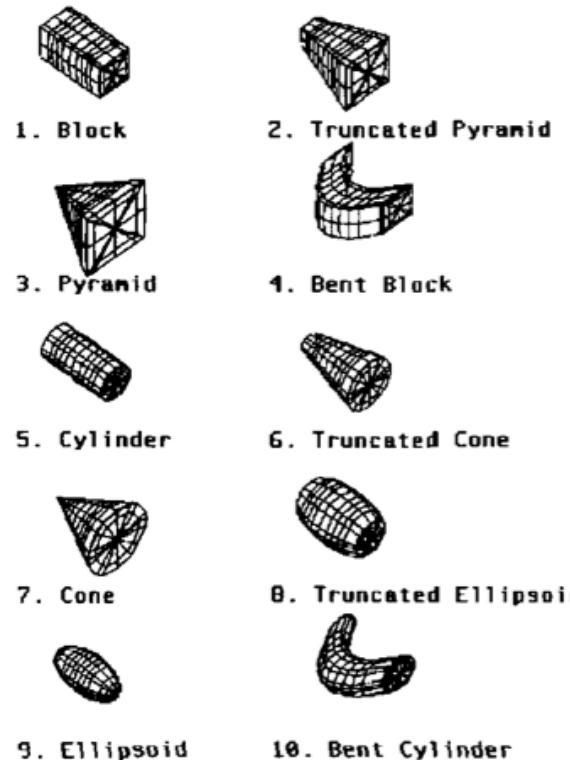
D. Marr

Fukushima's Neocognitron, 1983



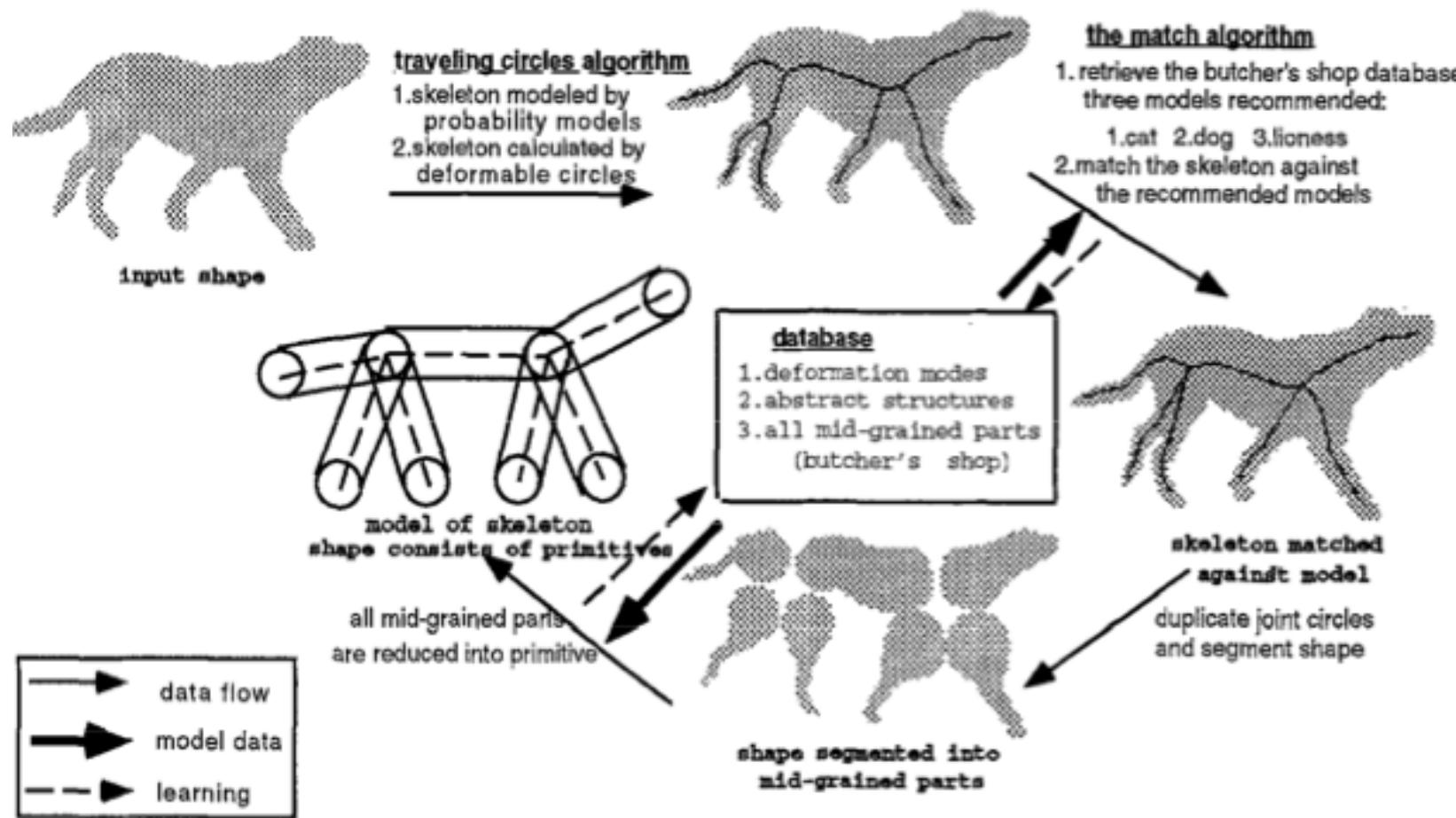
1	1	1	1	1	
2	/	/	/	/	
3	/	/	/	/	
4	/	/	/	/	
5	-	-	-	-	
6	1	1	1	1	
7	/	/	/	/	
8	/	/	/	/	
9	/	/	/	/	
10	-	-	-	-	
11	C	C	C	C	
12	C	C	C	C	
13	~	~	~	~	
14	?	?	?	?	
15	?	?	?	?	
16	?	?	?	?	
17	L	L	L	L	
18	L	L	L	L	
19	F	F	F	F	
20	7	7	7	7	
21	7	7	7	7	
22	Z	Z	Z	Z	
23	L	L	L	L	
24	L	L	L	L	
25	L	L	L	L	
26	+	+	+	+	
27	C	C	C	C	
28	C	C	C	C	
29	C	C	C	C	
30	C	C	-	C	
31	C	C	C	C	
32	~	~	~	~	
33	L	L	L	L	
34	Y	Y	Y	Y	
35	T	T	T	T	
36	A	A	A	A	
37	-	-	-	-	
38	4	4	4	4	

Perceptual Organization



- I. Biederman, Recognition-by-components: a theory of human image understanding, 1987
 S. Dickinson et. al., 3-D Shape Recovery Using Distributed Aspect Matching, 1992

Perceptual Organization



Generalized A* Parsing

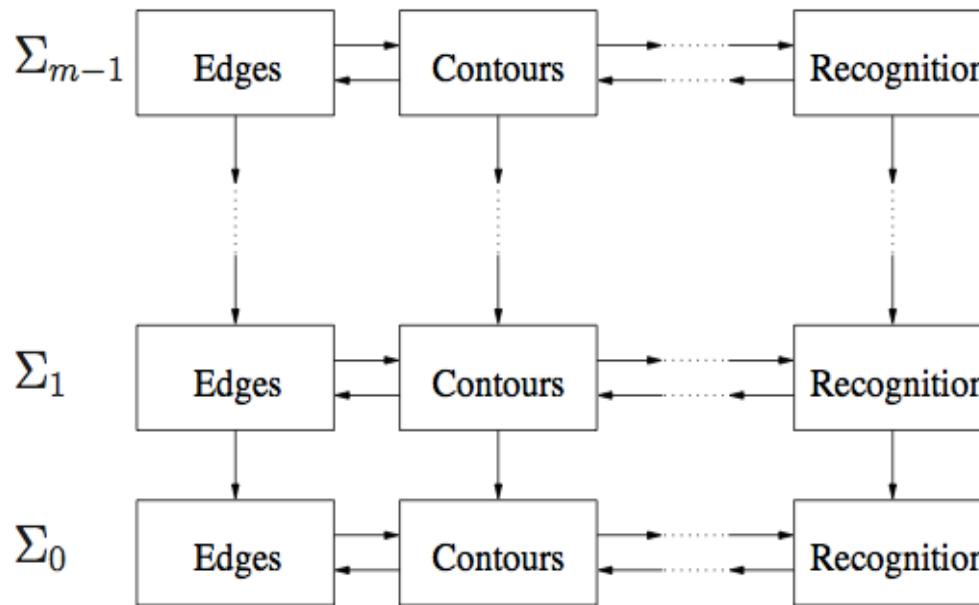
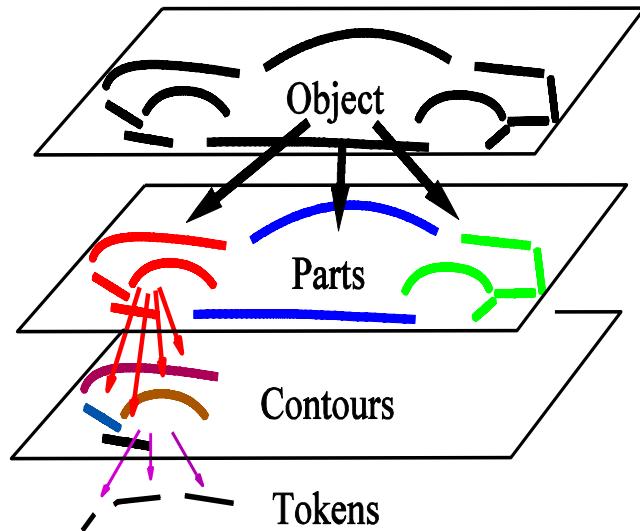
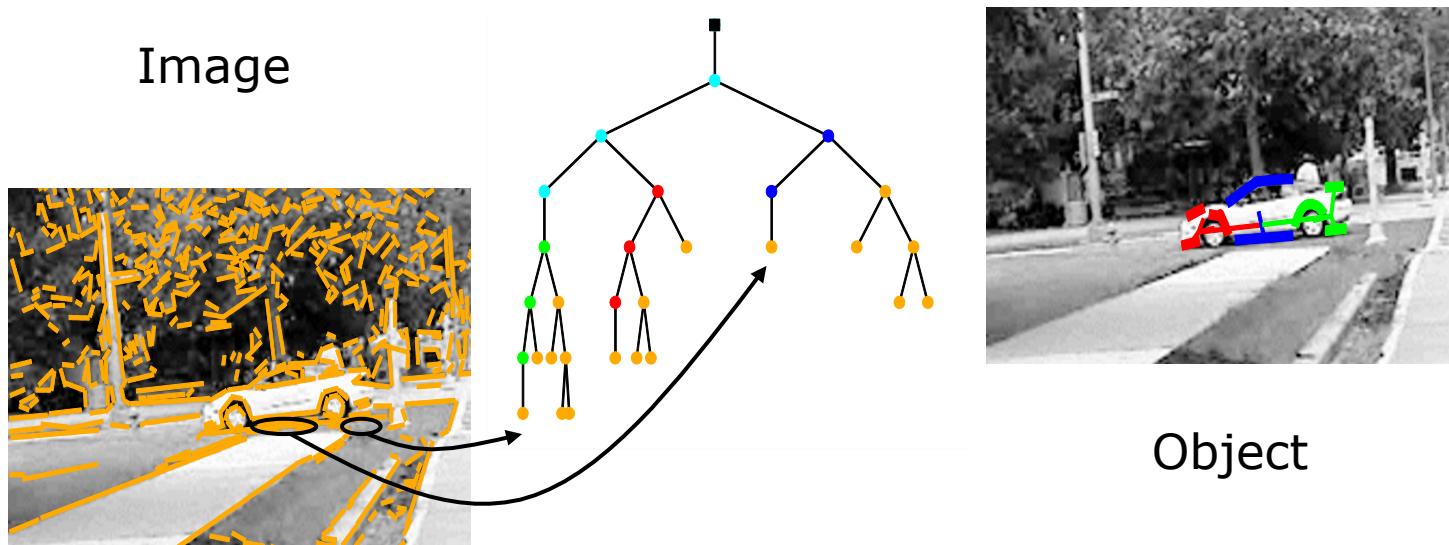


Figure 8: A vision system with several levels of processing. Forward arrows represent the normal flow of information from one stage of processing to the next. Backward arrows represent the computation of contexts. Downward arrows represent the influence of contexts.

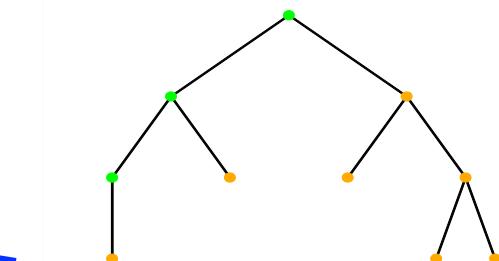
Hierarchical Shape Models



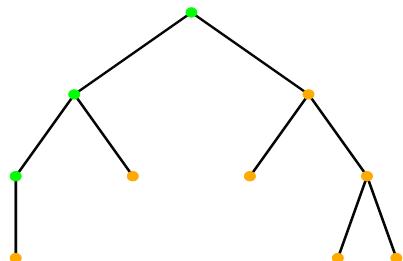
$$\begin{aligned}
 & \log P(\mathcal{T}_{1:C} | \mathcal{S}_1, \mathcal{P}_{1:3}) \\
 = & \log P_O(\mathbf{p}_O) + \sum_{p \in P} \log P_{p,O}(\mathbf{p}_p|_O) \\
 & + \sum_{p \in P} \sum_{c \in S_p} \log P_{c,p}(\mathbf{p}_c|_p) \\
 & + \sum_{c \in C} \log P(\mathcal{T}_{S_c} | \mathbf{p}_c) + c(\mathcal{T}_M)
 \end{aligned}$$



Composition of 'back'



Composition of 'back'

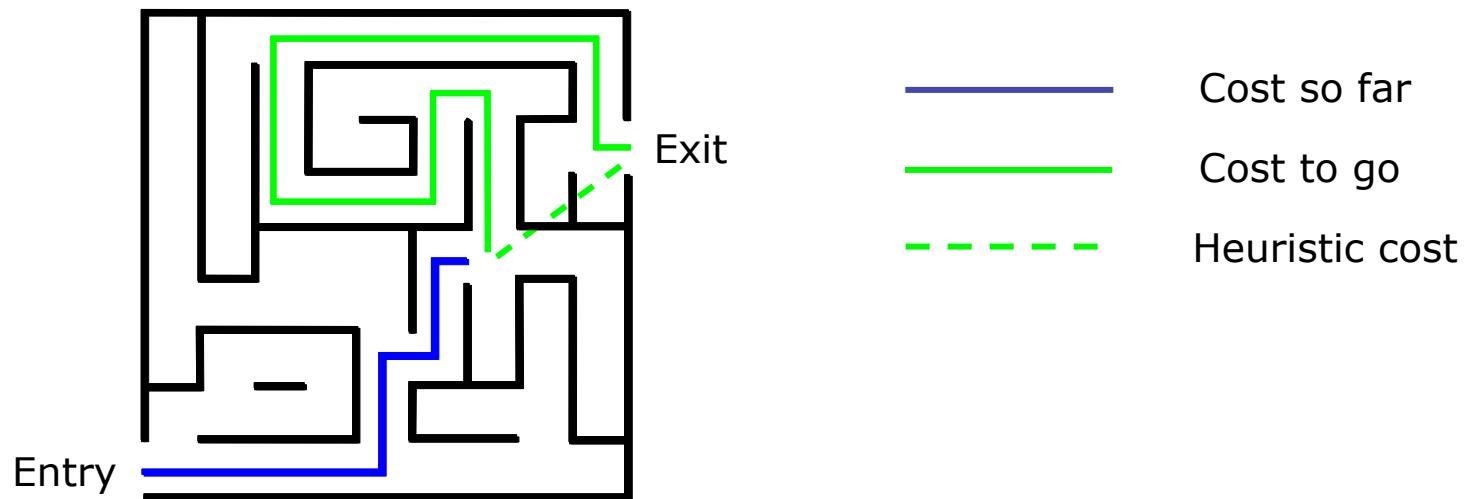


Problem: Too many options!

A* search & parsing

Dijkstra: prioritize based on ‘cost so far’ $c(s)$ $\{s : c(s) < c(S)\}$

A* : approximate ‘cost to go’ with heuristic $h(s)$ $\{s : c(s) + h(s) < c(S)\}$
prioritize using $c(s) + h(s)$



A* Parsing:

D. Klein and C. Manning, A* Parsing: Fast Exact Viterbi Parse Selection, (NAACL) 2003
P. Felzenszwalb, D. McAllester: The Generalized A* Architecture. (JAIR), 2007.

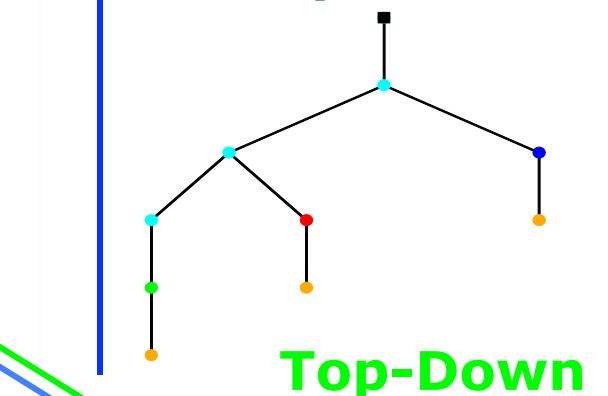
Our work: Object Parsing

Coarse-level parsing

Heuristics to Fine Level



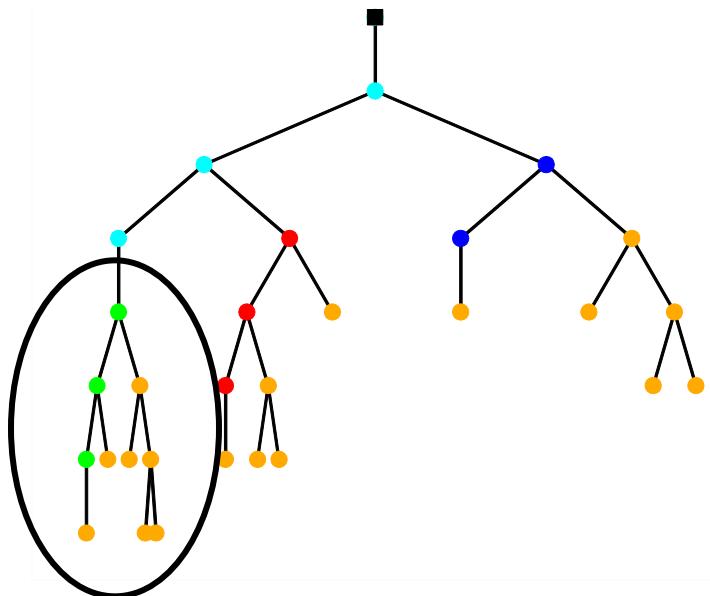
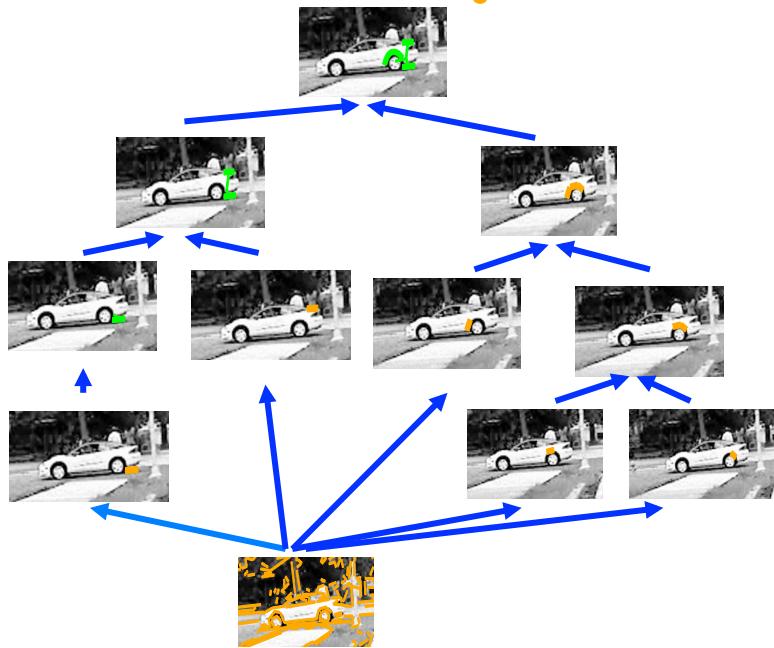
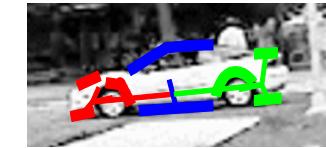
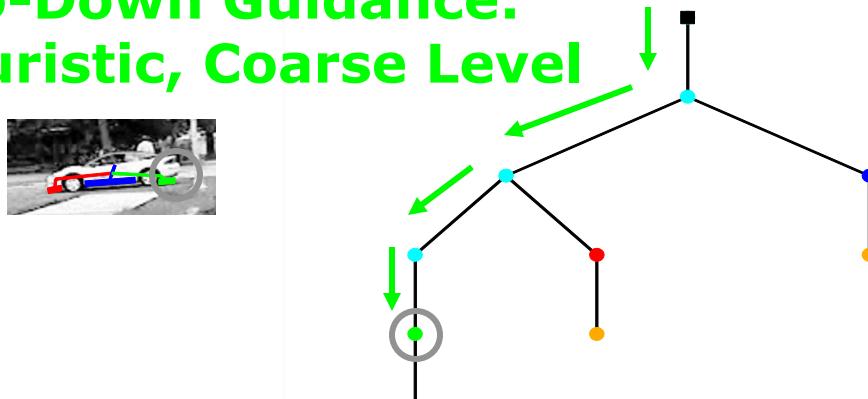
Bottom-Up



Top-Down

Fine-level parsing

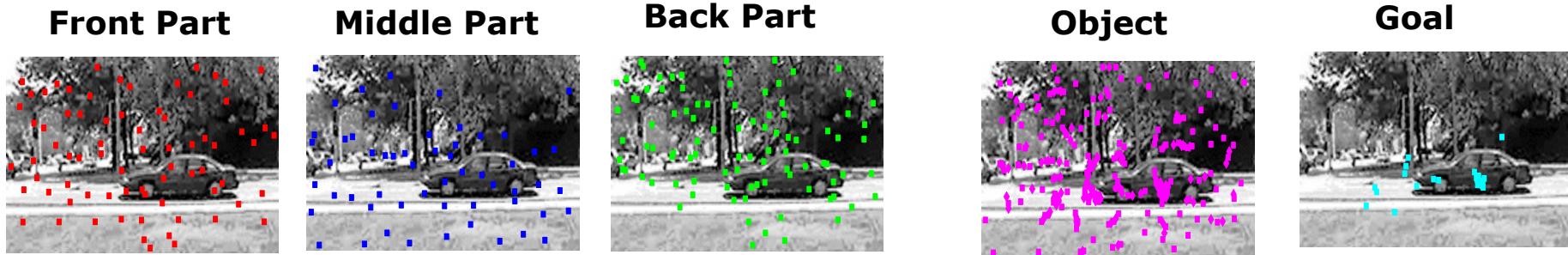
**Top-Down Guidance:
Heuristic, Coarse Level**



Bottom-Up Composition, Fine level

A* vs Knuth's Lightest Derivation (DP)

- A* Parsing

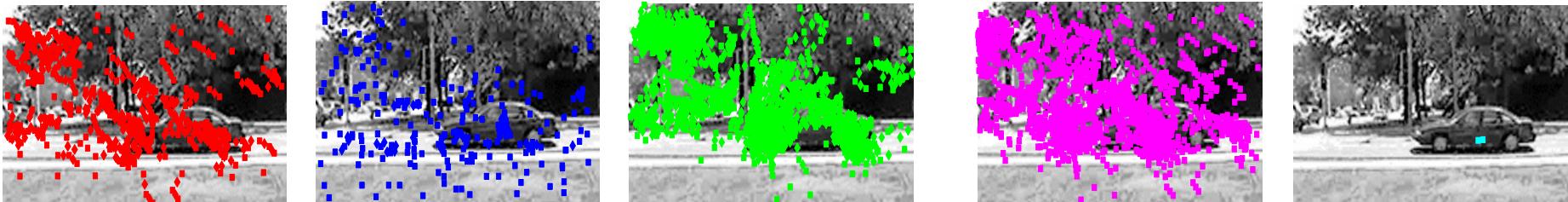


Coarse Level

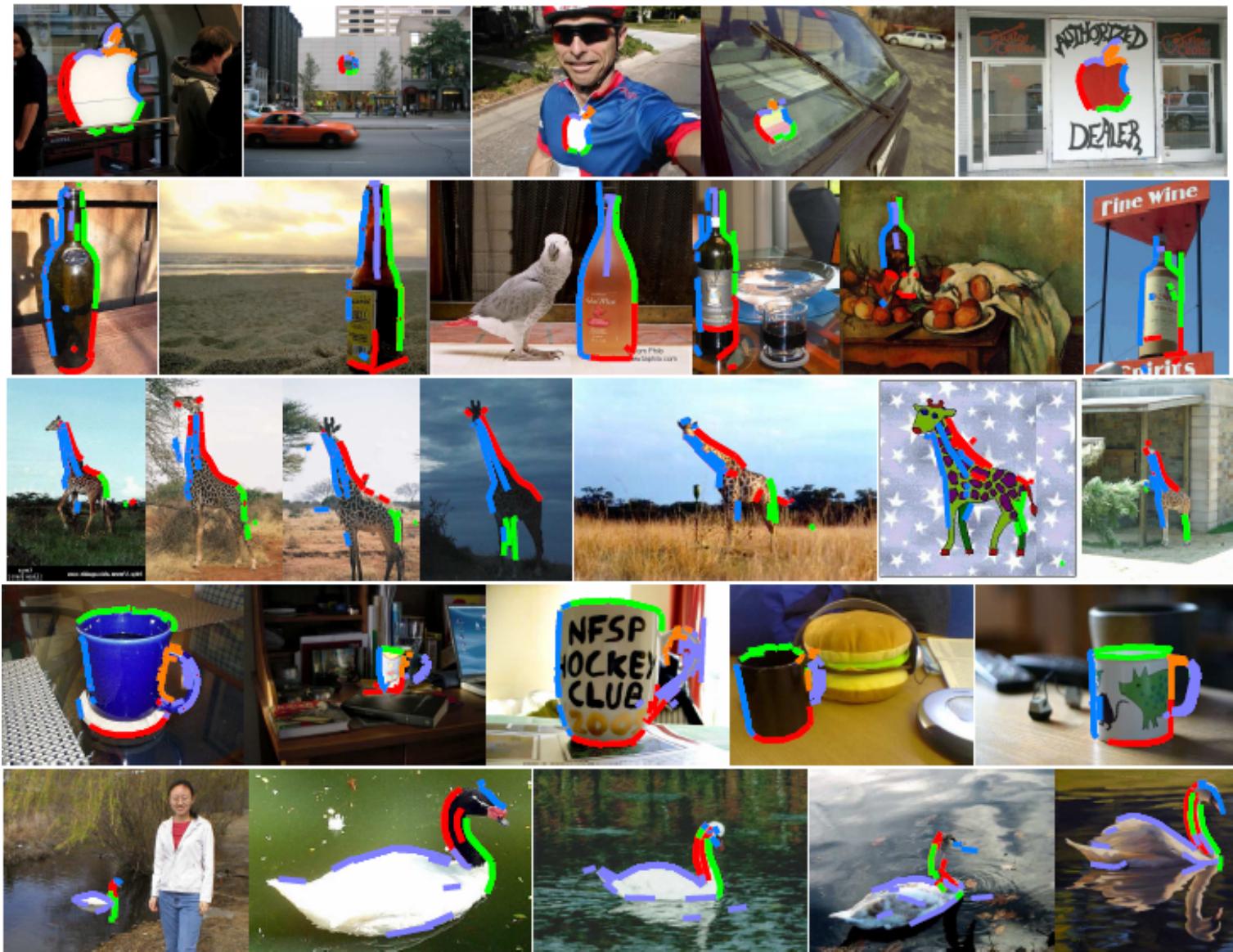


Fine Level

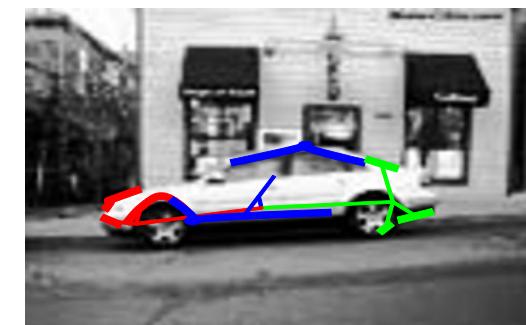
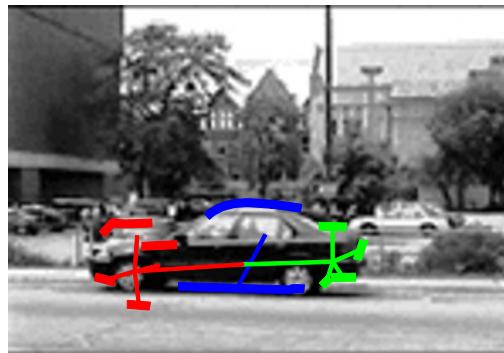
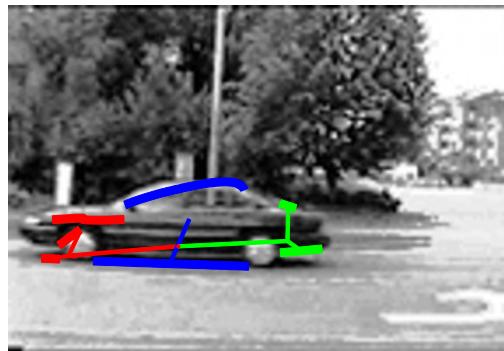
- KLD Parsing (only fine level)



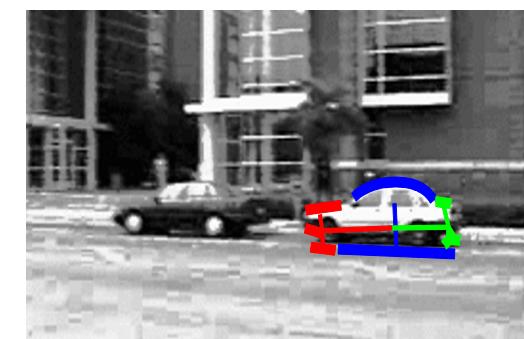
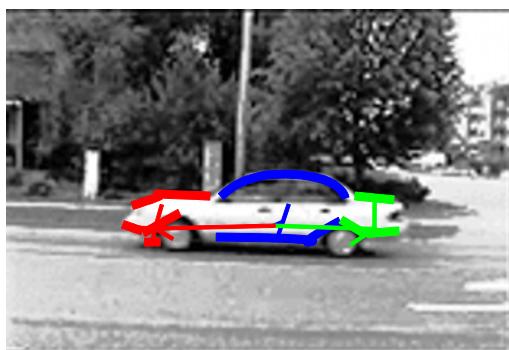
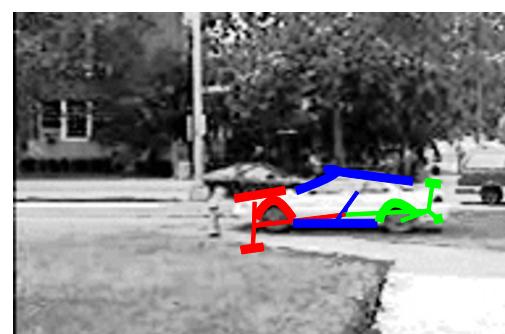
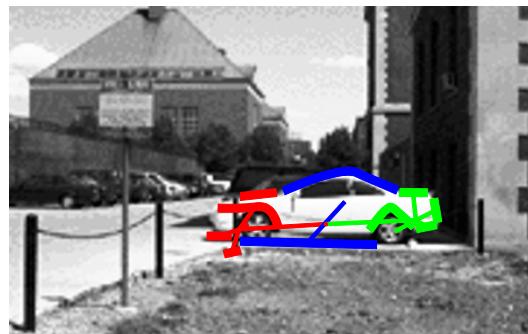
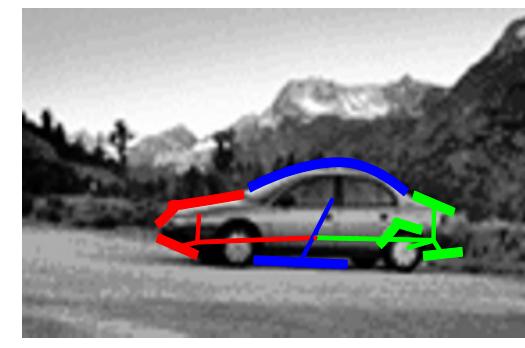
Parsing and localization results



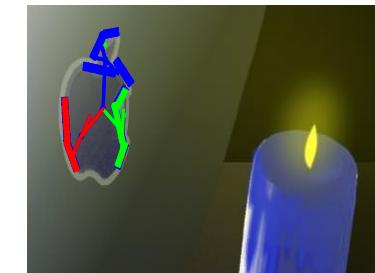
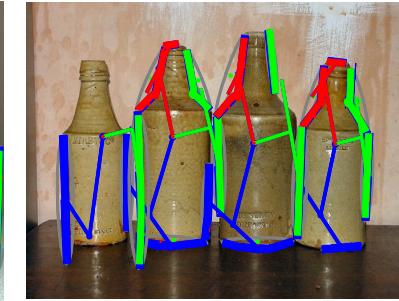
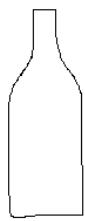
Parsing & Localization Results - I



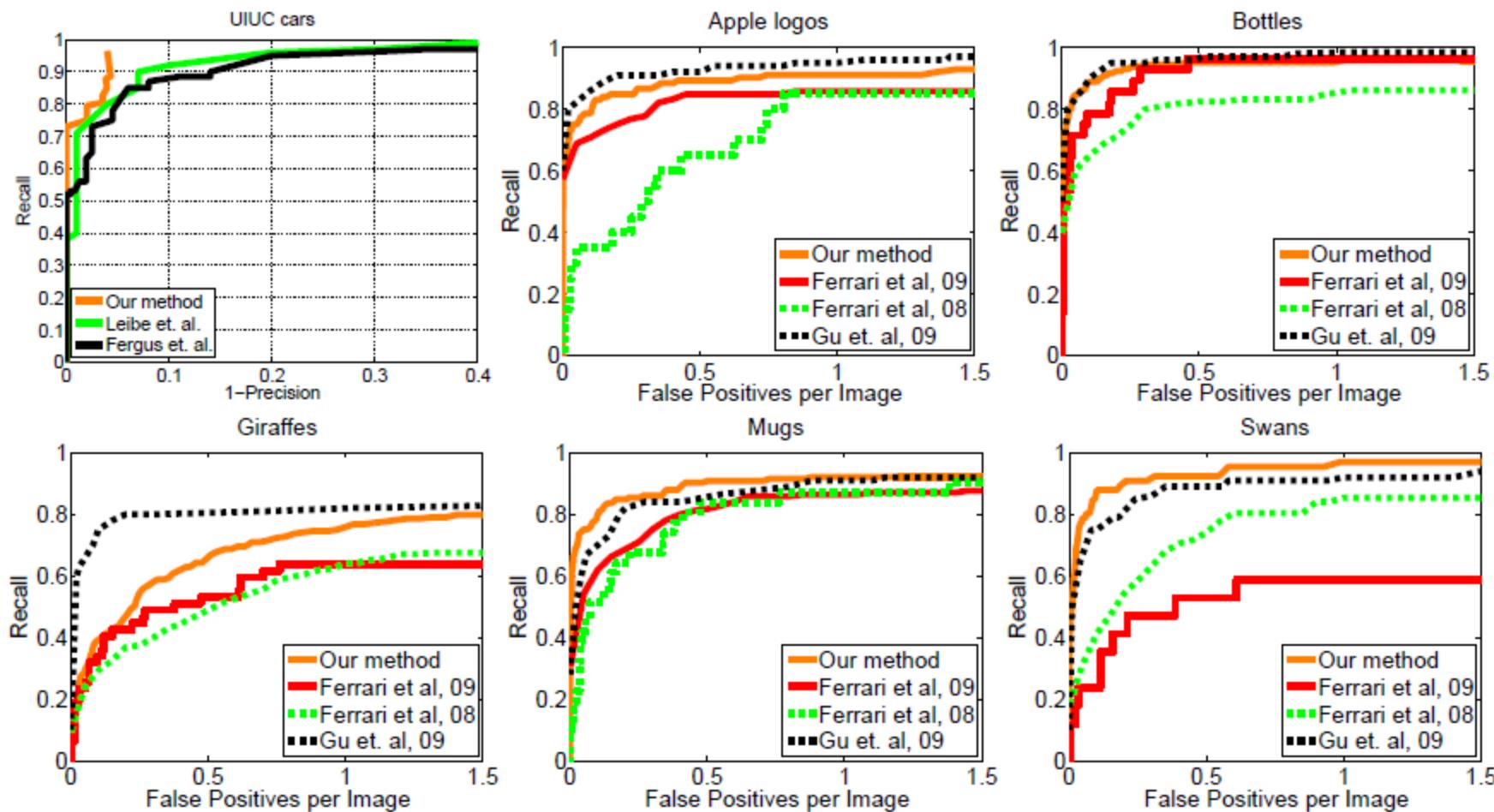
Parsing & Localization Results - II



Parsing & Localization Results - III

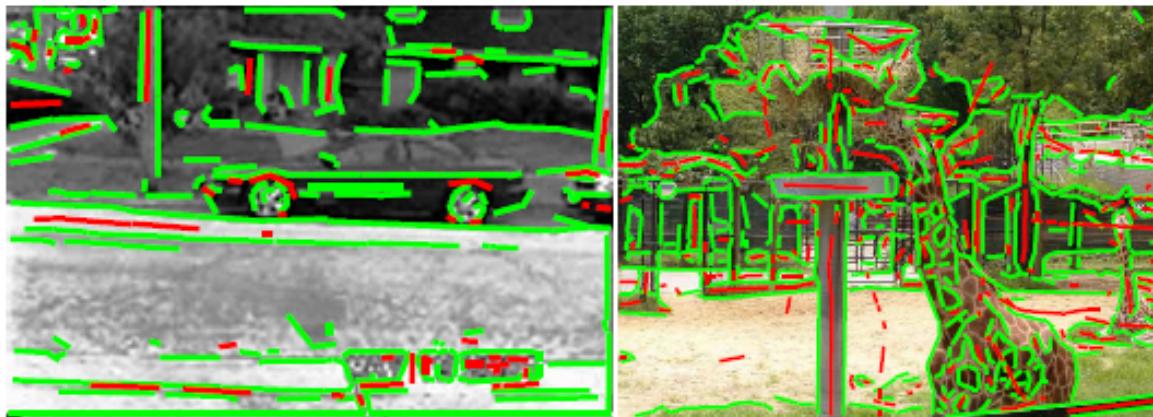


Benchmark results



Failure cases

Front-end failures



Missing appearance information/poor shape model

