

Learning Globally-Consistent Local Distance Functions for Shape-Based Image Retrieval and Classification

Andrea Frome
Yoram Singer
Fei Sha
Jitendra Malik

ICCV 2007



“dalmatian”



“dalmatian”



“buddha”

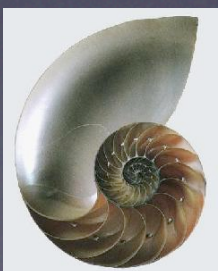


“dalmatian”

test image
nearest neighbor



“nautilus”



“beaver”



e.g., SIFT

feature comp

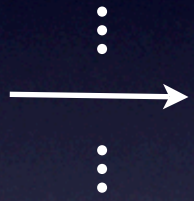
feature comp

distance/
similarity
function

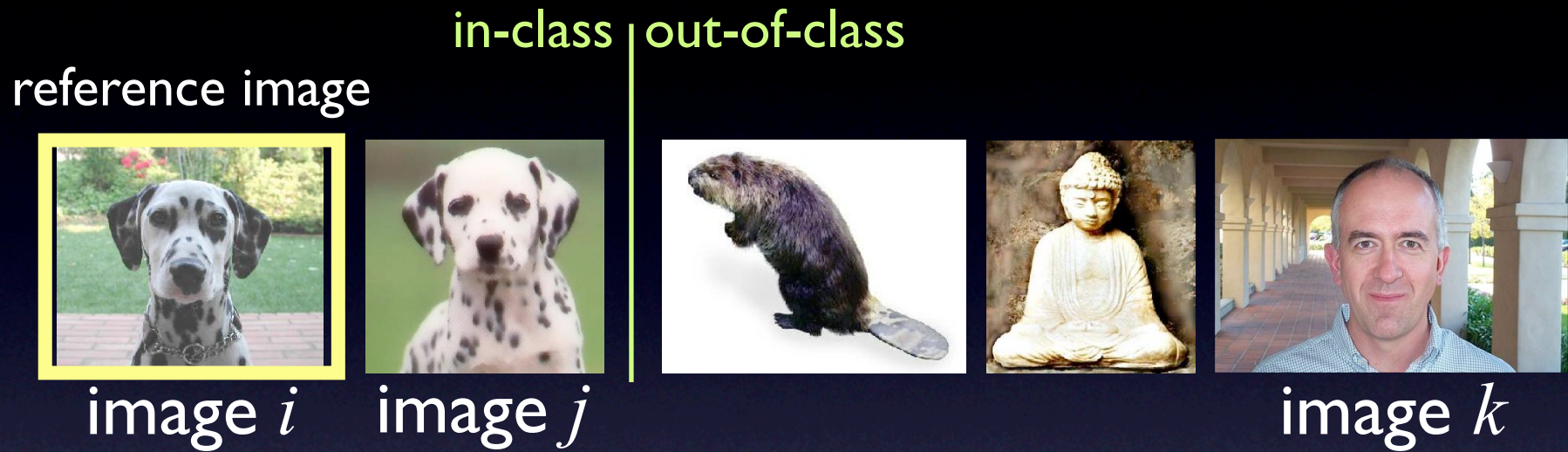
Learning

e.g., NN, SVM

learning
algorithm



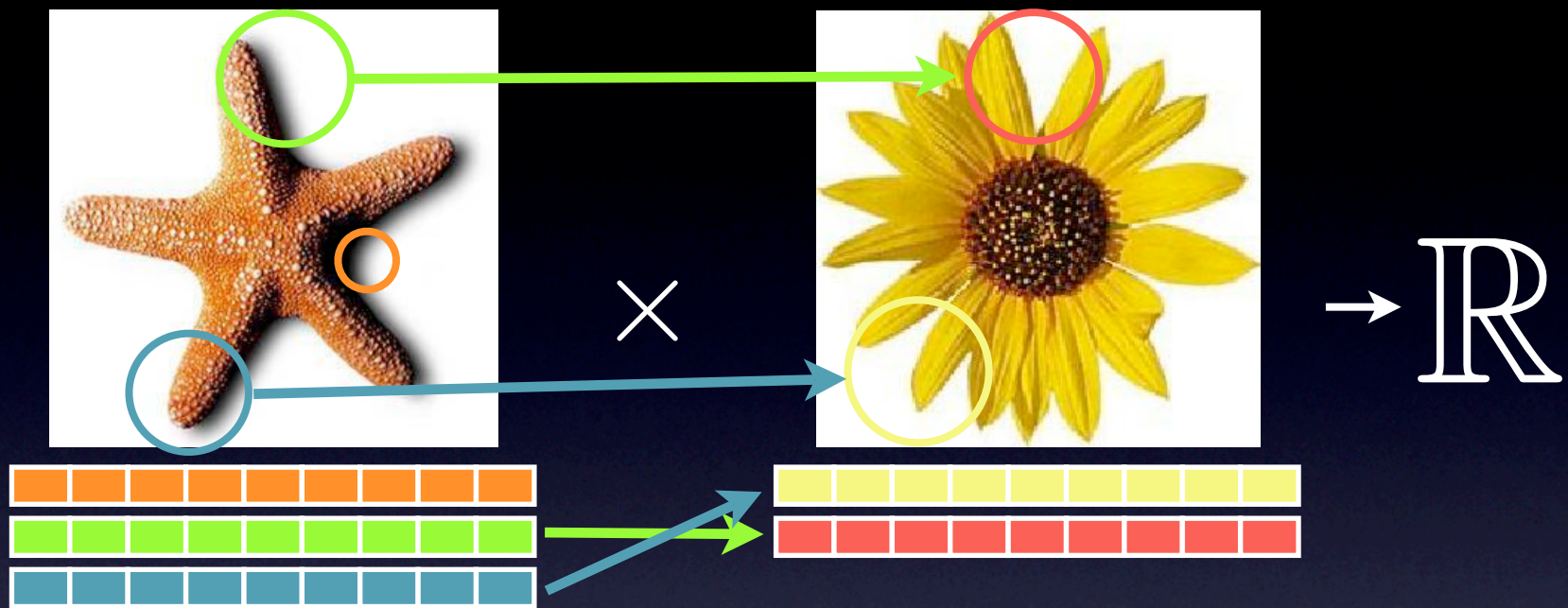
ranking: learn from **triplets** of training images



$$D(\text{image } k, \text{image } i) > D(\text{image } j, \text{image } i)$$

$$D_{ki} > D_{ji}$$

patch-based features



relaxations to correspondence

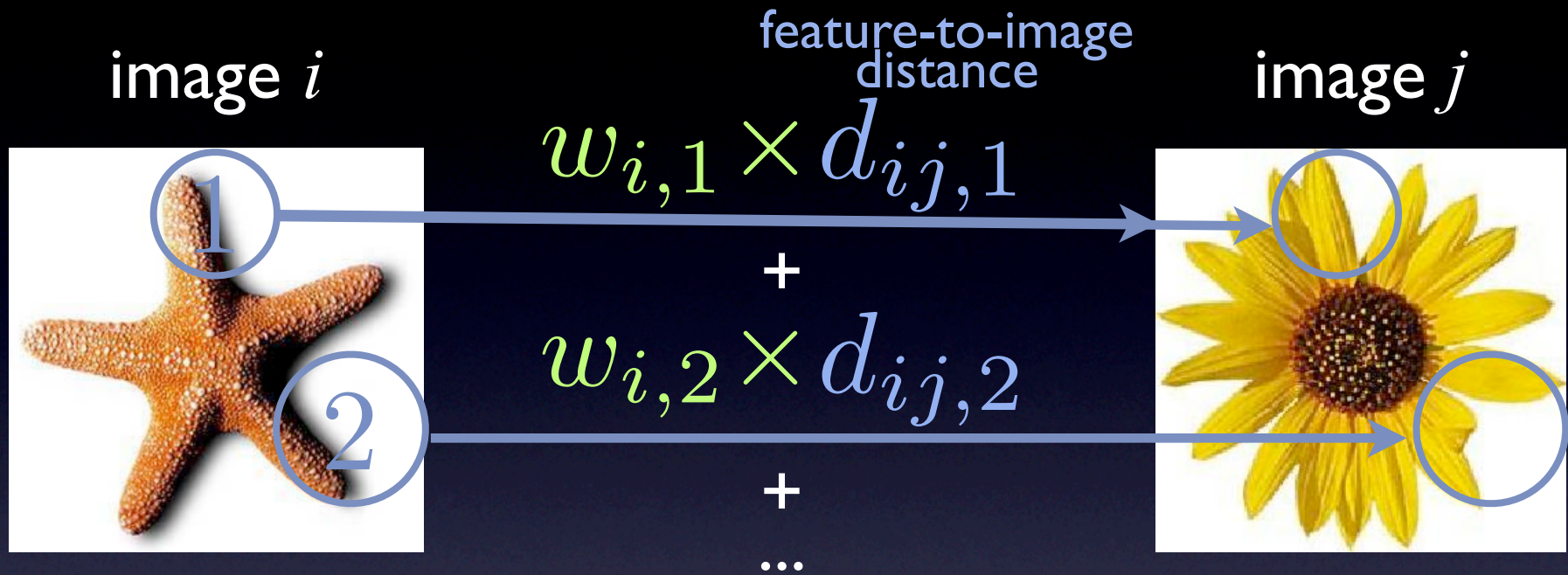
fast set matching (Grauman&Darrell 2006, Lazebnik, et al. 2006, Bosch, et al. 2007)

quantize feature space (bag of features) (Lazebnik, et al. 2006)

ignore spatial information (Grauman&Darrell 2006)

use absolute position information (Zhang, et al. 2006, Lazebnik, et al. 2006, Mutch&Lowe 2006)

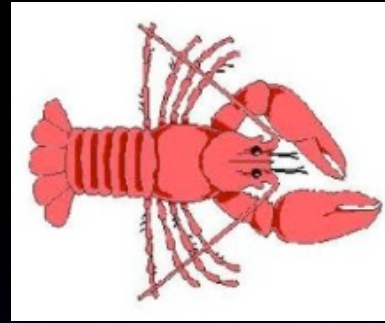
D_{ij} : distance from image i to image j
(not symmetric)



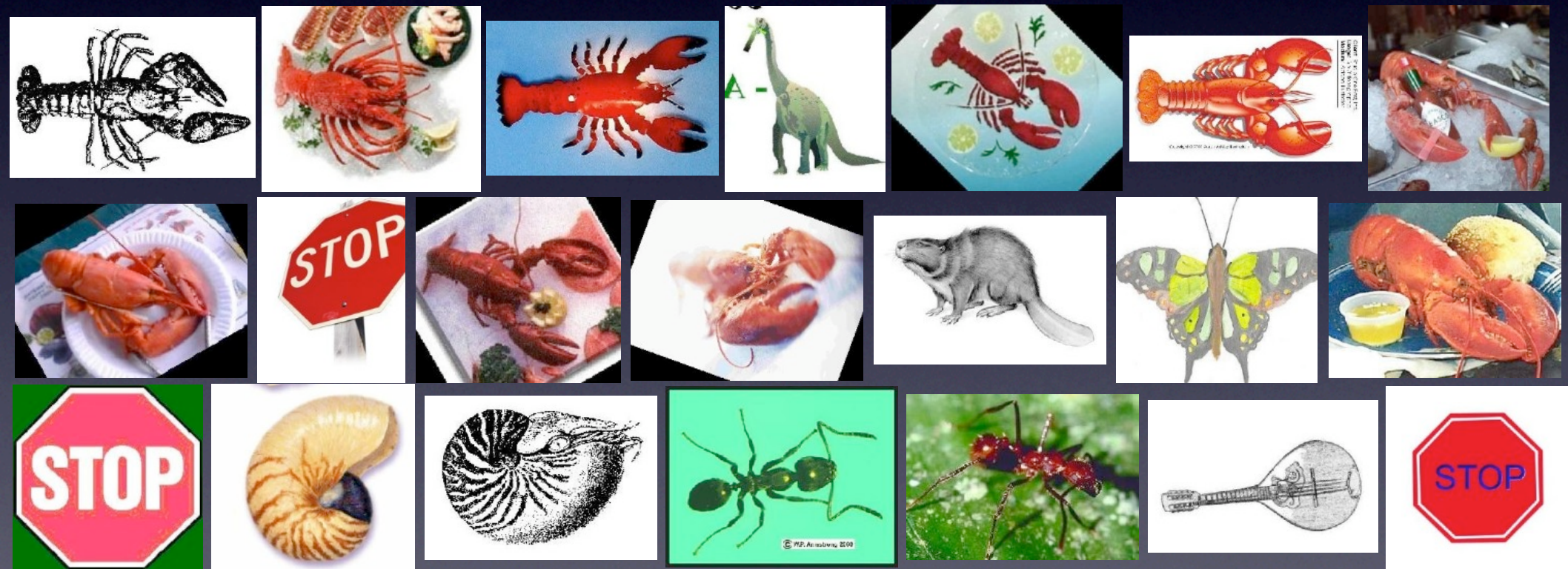
$$D_{ij} = \sum_{m=1}^M w_{i,m} d_{ij,m} = \mathbf{w}_i \cdot \mathbf{d}_{ij}$$

distance function can be evaluated from **image i** to
any other image

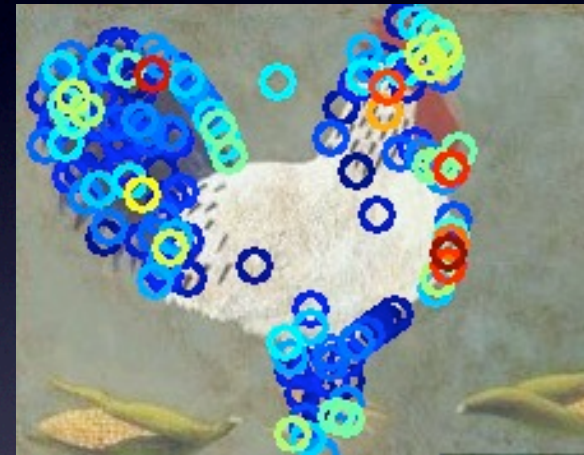
query image



retrieval results (weighted training images):



highest weight



lowest weight

why learn for every image?

clutter & occlusion



importance of a feature changes within a category

pose & articulation



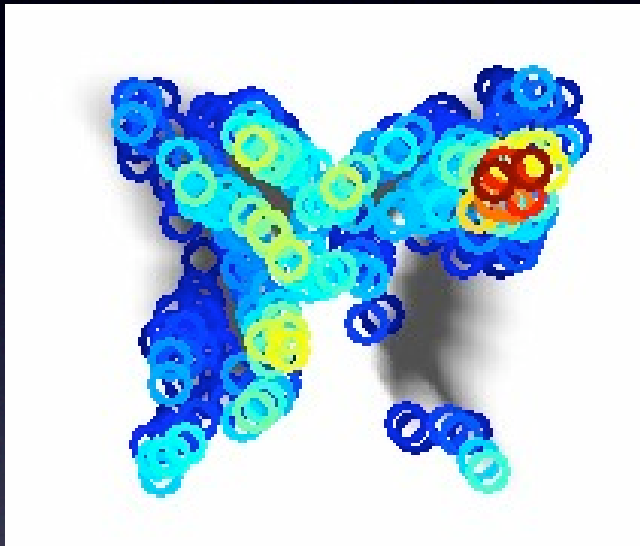
large variation



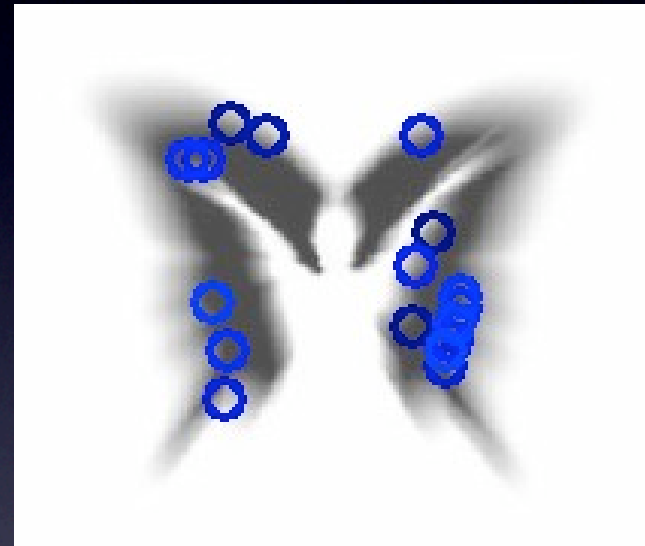
psychology: Rosch's family resemblances

Combine features in a way appropriate to *each image*

large-extent shape feature
(geometric blur)



color



- mathematical formulation
- relationship to other distance learning approaches
- selection of triplets
- results

“reference image”

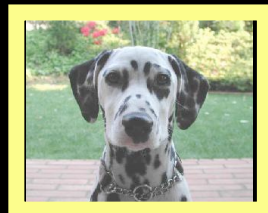


image i

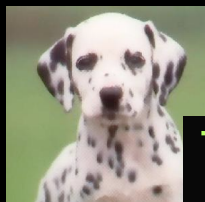


image j

w_j



image k

w_k



$$D_{ki} > D_{ji}$$

$$w_k \cdot d_{ki} > w_j \cdot d_{ji}$$

$$w_k \cdot d_{ki} - w_j \cdot d_{ji} > 0$$

w_k

w_j



W

d_{ki}

0

$-d_{ji}$

0



X_{ijk}

$$W \cdot X_{ijk} > 0$$

$$\mathbf{W} \cdot \mathbf{X}_{ijk} > 0$$

$$\mathbf{W} \cdot \mathbf{X}_{ijk} \geq 1$$

empirical loss: $\sum_{i,j,k \in \text{triplets}} [1 - \mathbf{W} \cdot \mathbf{X}_{ijk}]_+$

$$\min_{\mathbf{W}, \xi} \frac{1}{2} \|\mathbf{W}\|^2 + C \sum_{ijk} \xi_{ijk}$$

s.t. $\mathbf{W} \cdot \mathbf{X}_{ijk} \geq 1 - \xi_{ijk}$

$$\xi_{ijk} \geq 0$$

$$\mathbf{W} \succeq 0$$

Schultz, Joachims NIPS 2003
Frome, Singer, Malik NIPS 2006

problem scale

(15 images/category, 101 categories)

~1,200 features/image: weight vector has 1.8M elements
using in- vs. out-of-class,
exhaustive set of triplets is 31.8 M triplets

speeding it up

pare down to 15.7 M triplets

solve the dual problem

similar to on-line algorithms

early stopping: 10 hours to 1 hour

set trade-off parameter: one run through triplets

weight vectors are surprisingly sparse.

on average, 68% of weights are zero

Selecting triplets: 15 images/category

select 15.7 M out of 31.8 M triplets
many are easy, some are too hard

“reference”



easy triplet

“reference”



hard triplet

Heuristic using independent feature-to-image distances.

Relationship to other distance learning work

Zhang, Malik (CVPR 2003)
Bosch, Zisserman, Munoz (CVPR 2007)



learn a distance function for
all images
(global)

per category

one per image
(local)



Xing, Ng, Jordan, Russell (NIPS 2002)

Schultz, Joachims (NIPS 2003)

Shalev-Shwartz, Singer, Ng (ICML 2004)

Weinberger, Blitzer, Saul (NIPS 2005)

Globerson, Roweis (NIPS 2005)

Grangier, Monay, Bengio (ECML 2006)

Grauman, Darrell (NIPS 2006)

Varma, Ray (ICCV 2007)

Frome, Singer, Malik (NIPS 2006)

Frome, Singer, Sha, Malik (ICCV 2007)

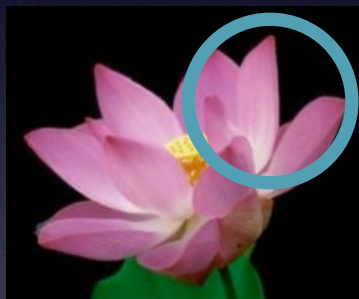
exploit collection of partial
descriptors
(patch-based features)

experiments

Caltech-101 (without using absolute position)

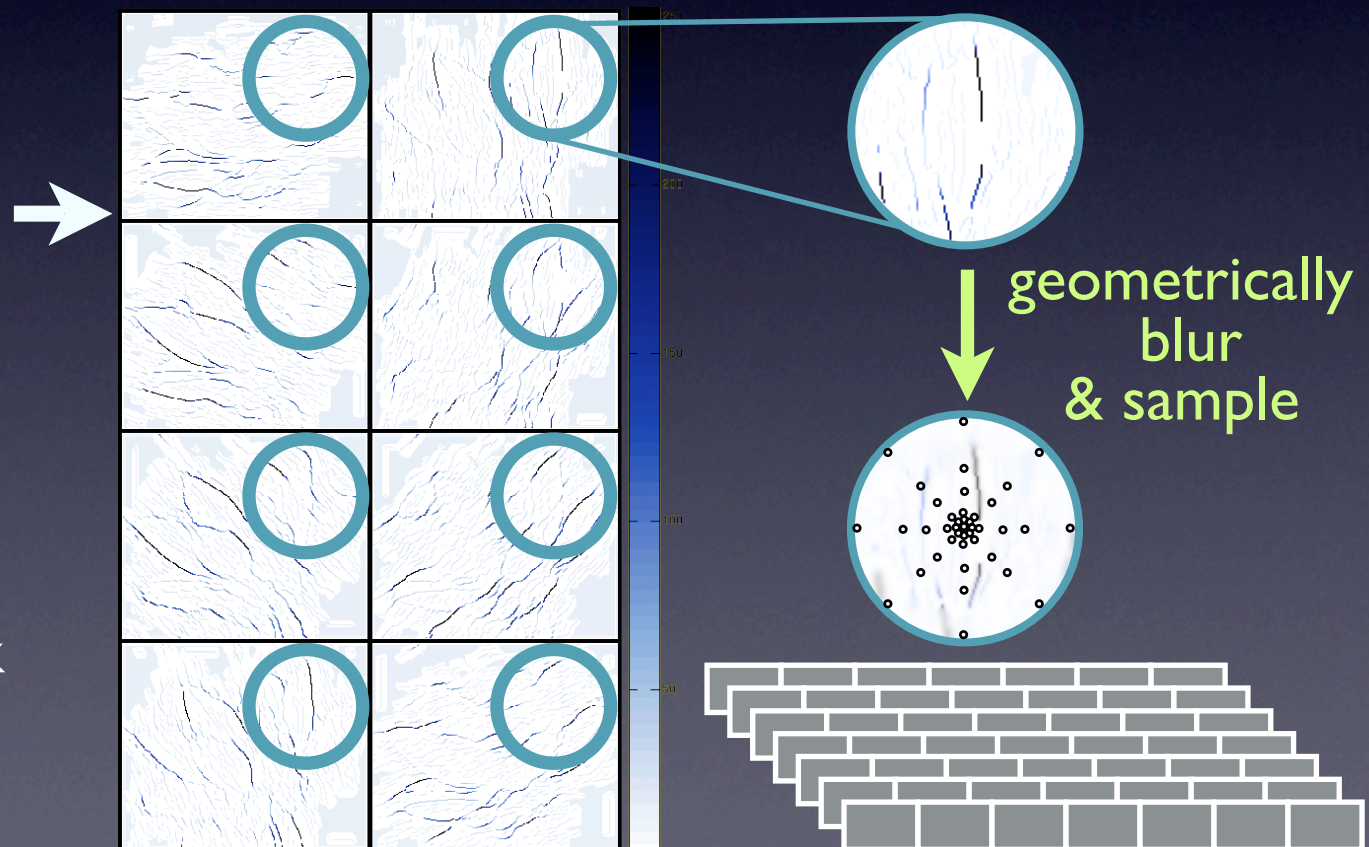
features: geometric blur (2 sizes) and color

L_2 feature-to-image distance



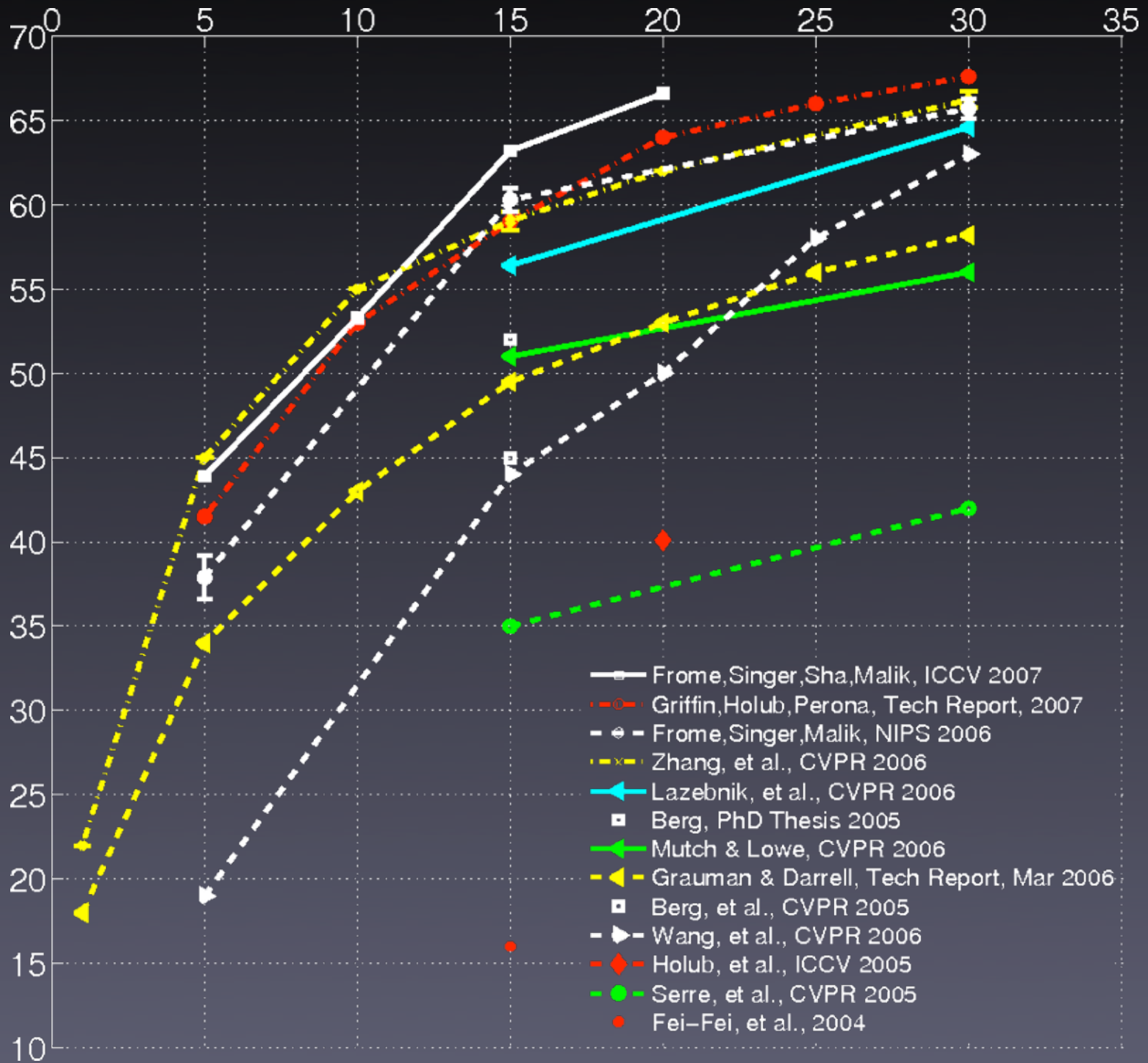
42 & 70 pixel
radius,
4 channels

Berg, Berg, Malik
CVPR 2005



training examples per class

mean recognition rate




<http://www.cs.berkeley.edu/~afrome/iccv2007>

Getting Started Latest Headlines goog mail dev handbook Machine Vision Algo... EB shuttle

Test image #3843











<< prev | next >> [All test images](#)



True classes:
butterfly

Predicted class: butterfly

fold #0
image #3843

 <p>13.170024 3825 butterfly</p>	 <p>13.536894 3817 butterfly</p>	 <p>14.719784 7635 rooster</p>	 <p>14.875127 3836 butterfly</p>	 <p>15.006959 3857 butterfly</p>	 <p>15.118133 3860 butterfly</p>	 <p>15.196454 8104 starfish</p>	 <p>15.319606 3610 brain</p>	 <p>15.400794 3820 butterfly</p>	 <p>15.405533 3889 butterfly</p>
---	---	--	---	--	---	--	---	---	---

thank you.