



Context-based Object Recognition - Putting context into vision

Phil Crosby and Zhe Lin

CMSC 828J

May 04, 2006



Outline

- n Introduction to context
- n Overview of Cognitive Science Literature
- n Matching words and images
- n Statistical contextual priming
- n Contextual priors for object recognition
- n Other contextual models
- n Summary
- n References



Definitions of Context

- n Any data or meta-data not directly produced by object of interest.
Includes:
 - n Neighbor-based context: Nearby image data
 - n Scene-based context: Scene information
 - n Object-based context: Presence, locations of other objects



Contextual Viewing of Visual Attention

Marvin Chun



Chun on Attention

- n Attention must prioritize relevant information in the face of information overload
- n Bold and flamboyant bottom-up cues can produce rapid search
- n Chun argues normal scenes have **too many** bottom-up cues to be useful



Chun on Attention...

- n Context dictates what should receive attention and what should be ignored
 - n Lights while driving
- n Context guides eye movements to the important parts of a scene, which are then fixated upon with foveal vision



Yarbus, A.L. (1967) *Eye Movements and Vision*, Plenum Press



Biederman - Context

- n If context guides our attention, it basically facilitates a more efficient search
- n Biederman's (1981) influential theory: schema representations for a scene specify the range of plausible objects that can occur, and their positioning relative to each other
- n Schemas acquired rapidly, with a glance



Potential Benefits of Context

- n Eliminate ambiguity!
- n Place constraints on the types of objects (Chun, Biederman)
 - n Reduces computational complexity when searching memory for an identity match for a new object



Potential Benefits of Context...

- n Focus attention, giving more perceptual processing to more relevant image areas
- n Exploit redundancy and invariance in a scene.
 - n Discard irrelevant or non-changing details
 - n Increase predictability to save time



Cognitive science studies on context



What letter is this?

A

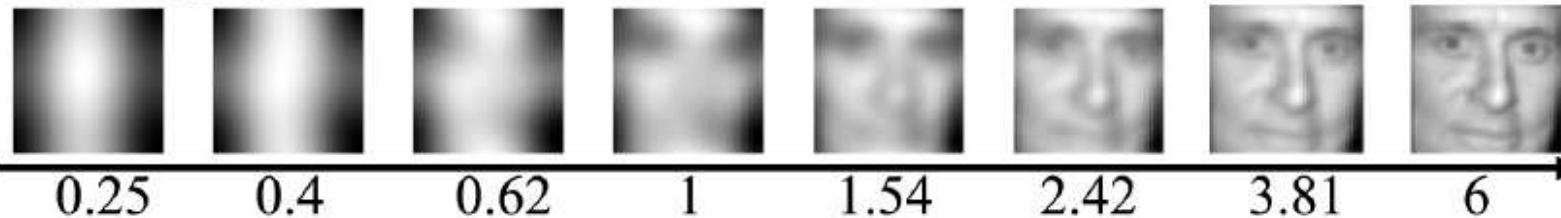


Selfridge - Disambiguation

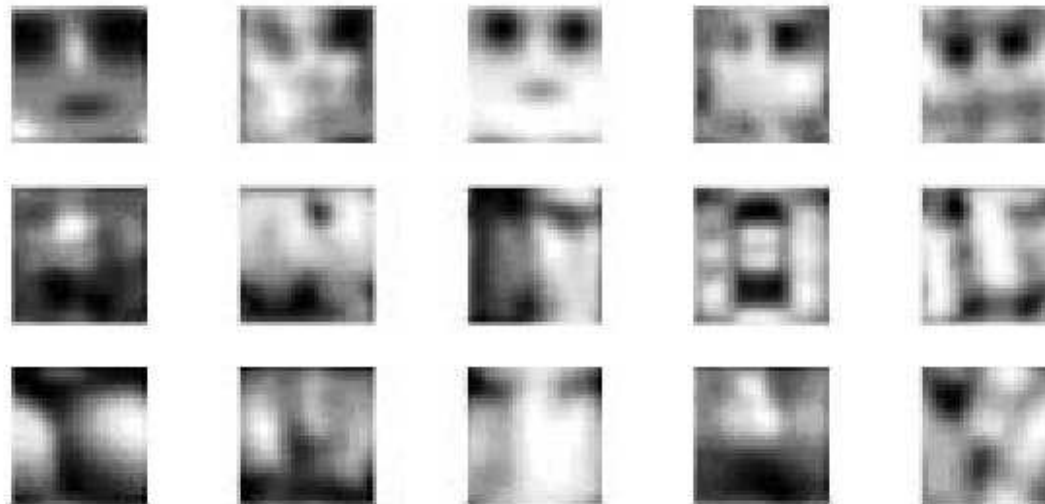
TAE CAT

Torralba – Face detection in impoverished images

Low resolution

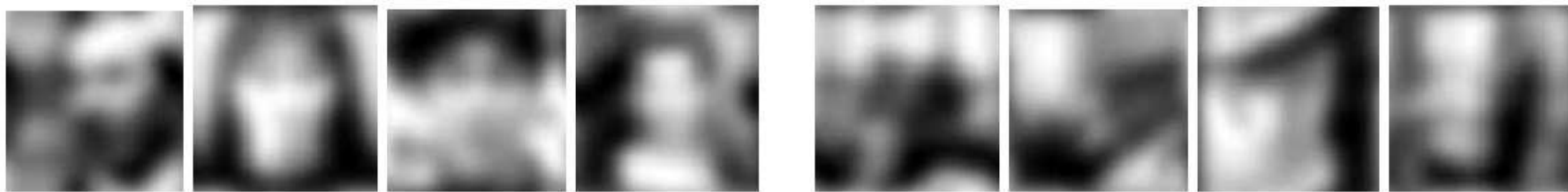
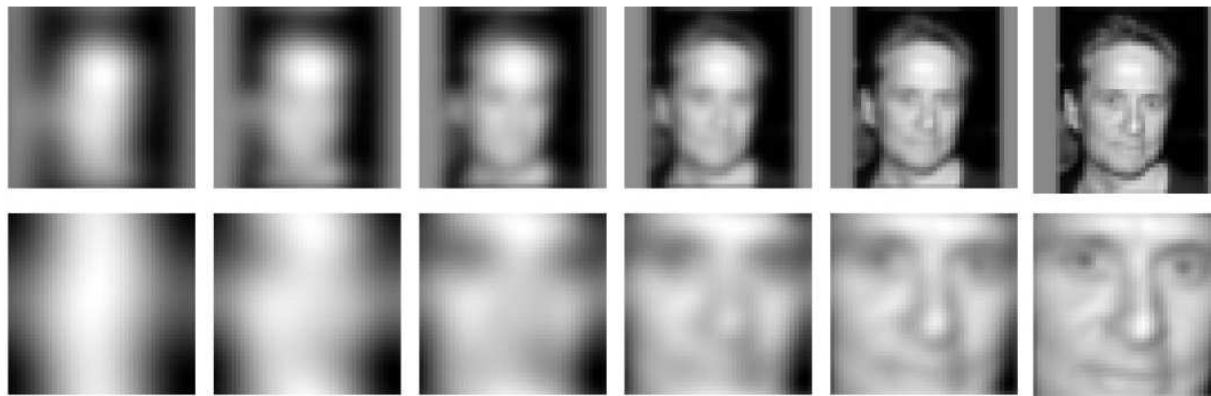


Versus distracter shapes:

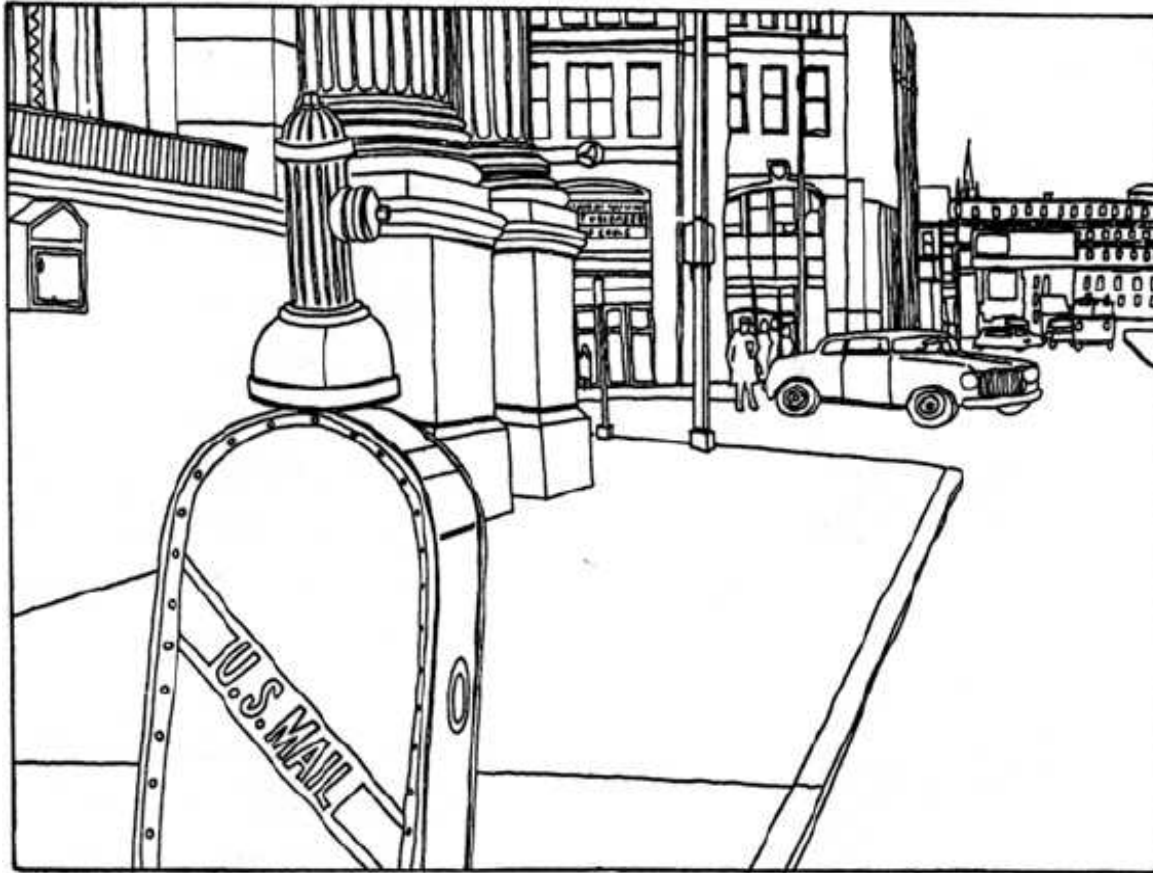


Torralba – Face detection in impoverished images

Same resolution, but trade details of face for contour of head



Where is the fire hydrant?

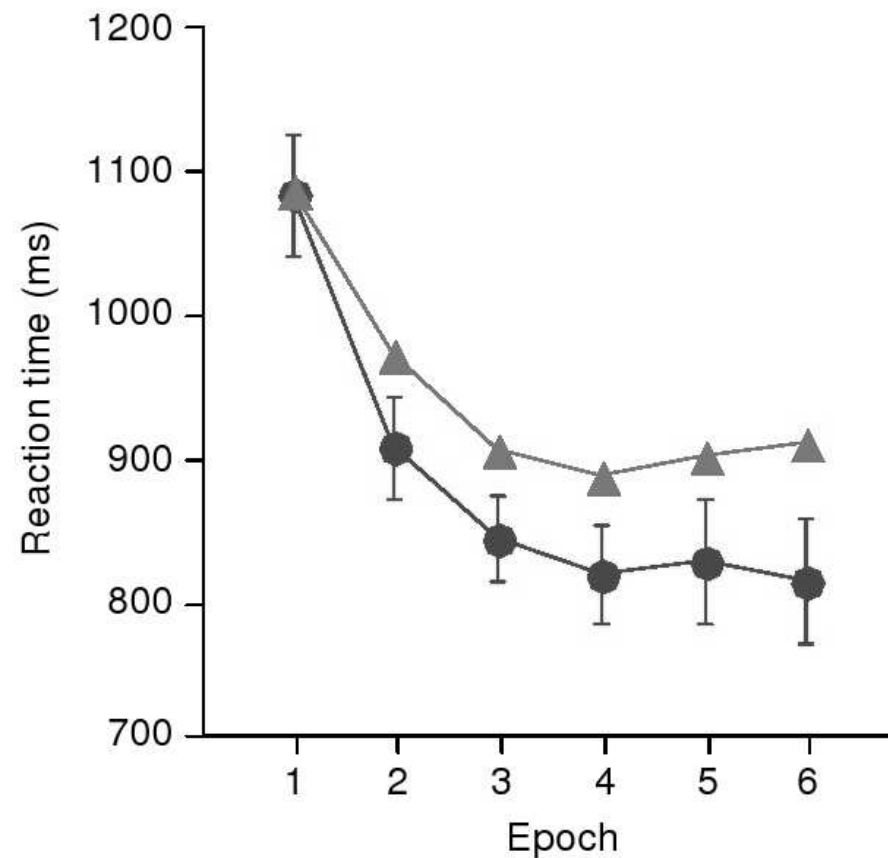
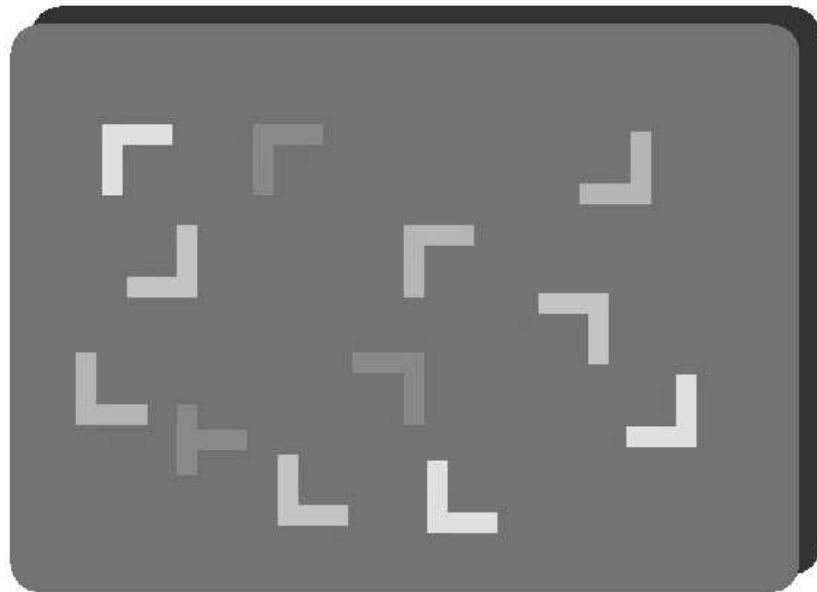


Biederman, I. et al. (1982) Scene perception: detecting and judging objects undergoing relational variations. *Cognit. Psychol.* 14, 143-177

I think I know that face...

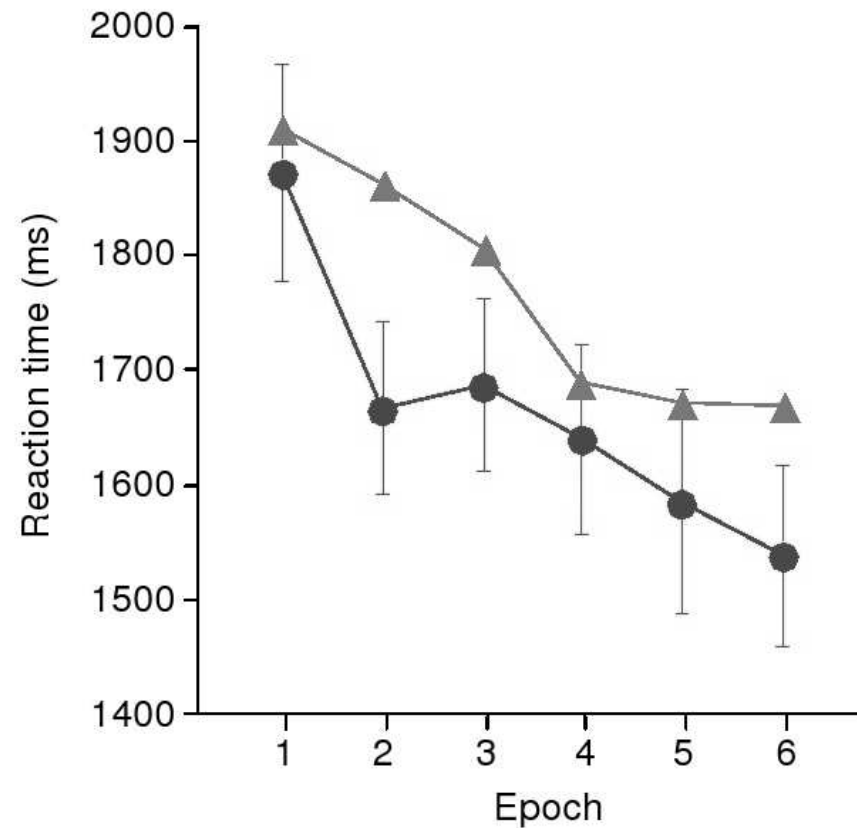
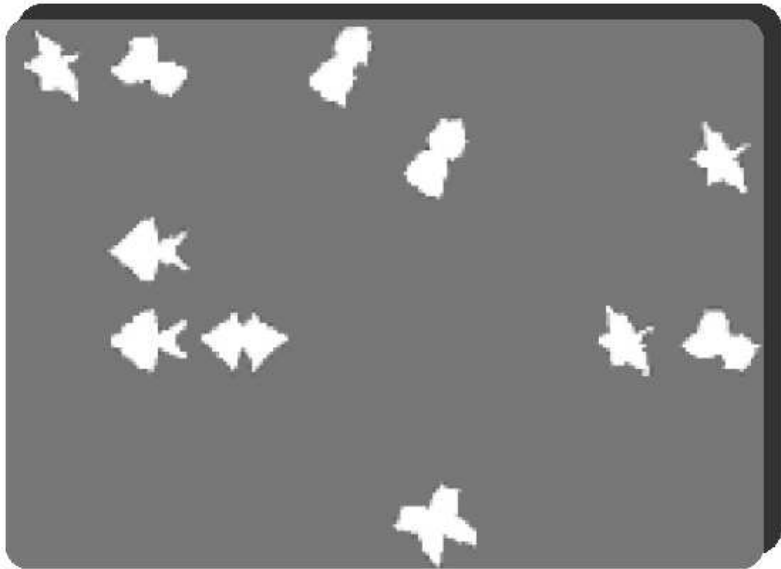


Spatial contextual cueing



Chun, M.M. and Jiang, Y. (1998) Contextual cueing: implicit learning and memory of visual context guides spatial attention. *Cognit. Psychol.* 36, 28–71

Object cueing



Chun, M.M. and Jiang, Y. (1999) Top-down attentional guidance based on implicit learning of visual covariation. *Psychol. Sci.* 10, 360–365



Henderson, Hollingworth – High Level Scene Perception

- n How do we explain the results of these studies?
- n Three hypotheses for human object identification in scenes:



Perceptual Schema Model (Biederman's view)

- n Expectations derived from knowledge about the composition of a scene type interact with the perceptual analysis of object tokens in that scene.
 - n Facilitates perceptual analysis of scene-consistent objects
 - n Inhibits analysis of scene-inconsistent objects



Priming Model

- n Recognition of a scene “primes” the memory associated with objects that should occur in that scene.
- n Criterion for matching a visual observation to an object in memory **is lower** for object memories that are consistent with that scene



Functional Isolation Model

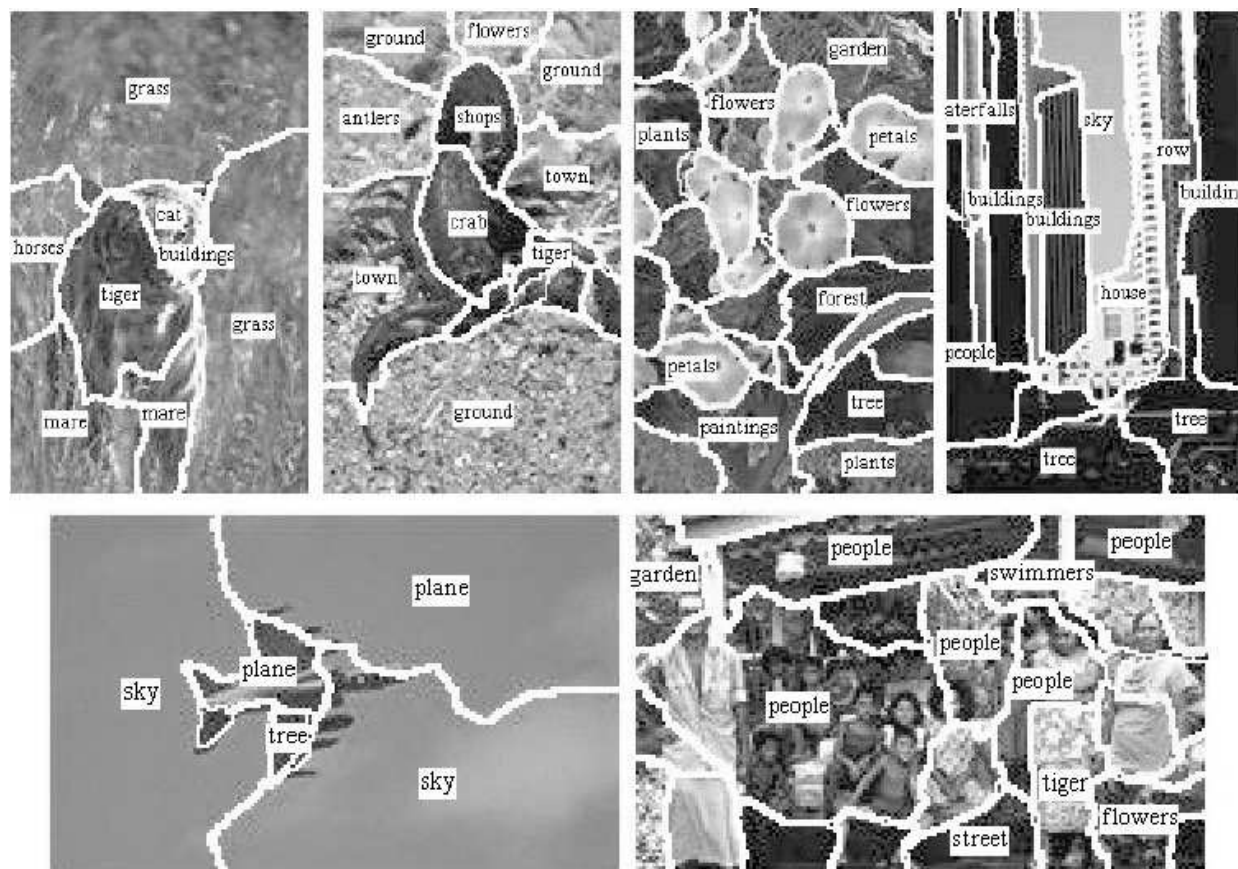
- n Object identification is isolated from expectations derived from scene knowledge



Which model is accurate?

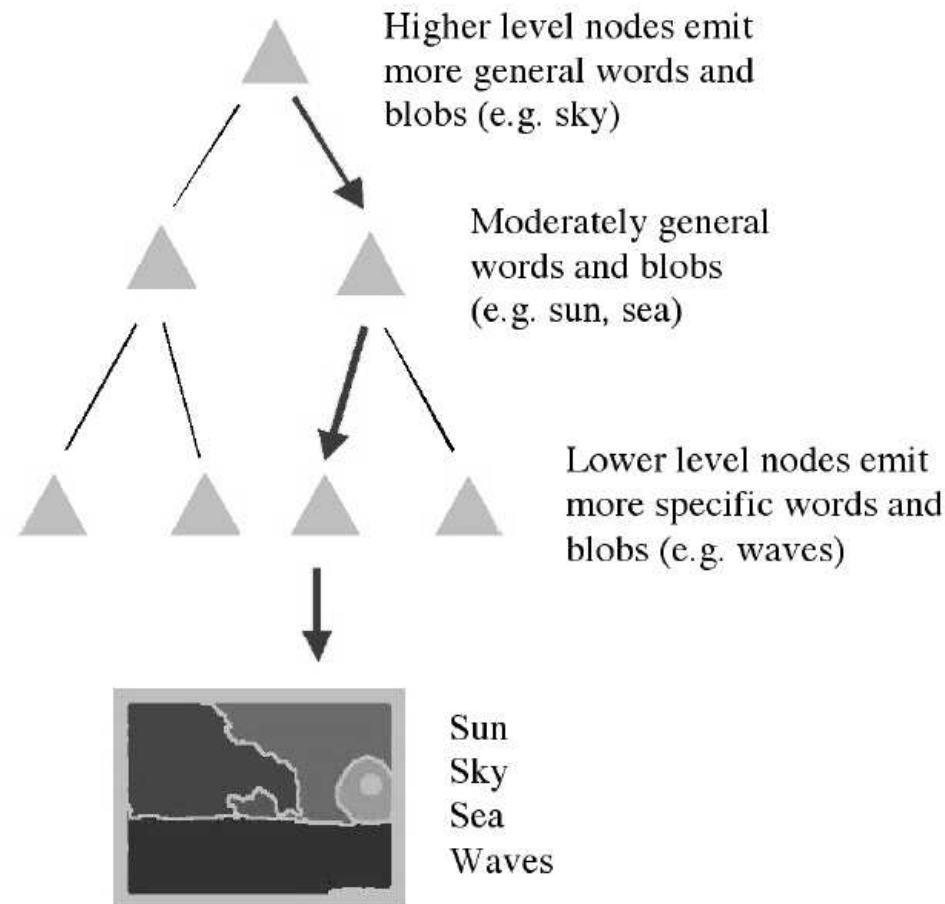
- n Who knows?
- n Henderson & Hollingworth seem to find that the functional isolation model best describes how humans do object recognition

Barnard et al – Matching words and Pictures



Barnard, K., Duygulu, P., Forsyth, D., de Freitas, N., Blei, D. M., and Jordan, M. I. 2003. Matching words and pictures. *J. Mach. Learn. Res.* 3 (March 2003), 1107-1135.

Matching words and Pictures - Hierarchical model for annotation



Barnard, K., Duygulu, P., Forsyth, D., de Freitas, N., Blei, D. M., and Jordan, M. I. 2003. Matching words and pictures. *J. Mach. Learn. Res.* 3 (March 2003), 1107-1135.

Matching words and Pictures - Hierarchical model for annotation

$$p(D|d) = \sum_c p(c) \prod_{w \in W} \left[\sum_l p(w|l, c) p(l|d) \right]^{\frac{N_w}{N_{w,d}}} \prod_{b \in B} \left[\sum_l p(b|l, c) p(l|d) \right]^{\frac{N_b}{N_{b,d}}}$$

$$\begin{aligned} p(w|B) &\propto \sum_c p(c) p(w|c) p(B|c) \\ &= \sum_c p(c) \left[\sum_l p(w|l, c) p(l|c) \right] \prod_{b \in B} \left[\sum_l p(b|l, c) p(l|c) \right]^{\frac{N_b}{N_{b,d}}} \end{aligned}$$

Matching words and Pictures - leveraging context

- n Leverage clusters in the hierarchical model to help make clustering more accurate. Clusters are the context.
- n **Approach 1** – linking word emission and region emission probabilities with **mixture weights**:

$$p(D|d) = \sum_c p(c) \prod_{w \in W} \left[\sum_l p(w|l, c) p(l|B, c, d) \right]^{\frac{N_w}{N_{w,d}}} \prod_{b \in B} \left[\sum_l p(b|l, c) p(l|d) \right]^{\frac{N_b}{N_{b,d}}},$$

where we stipulate that

$$p(l|B, c, d) \propto \sum_{b \in B} p(l|b, c, d).$$

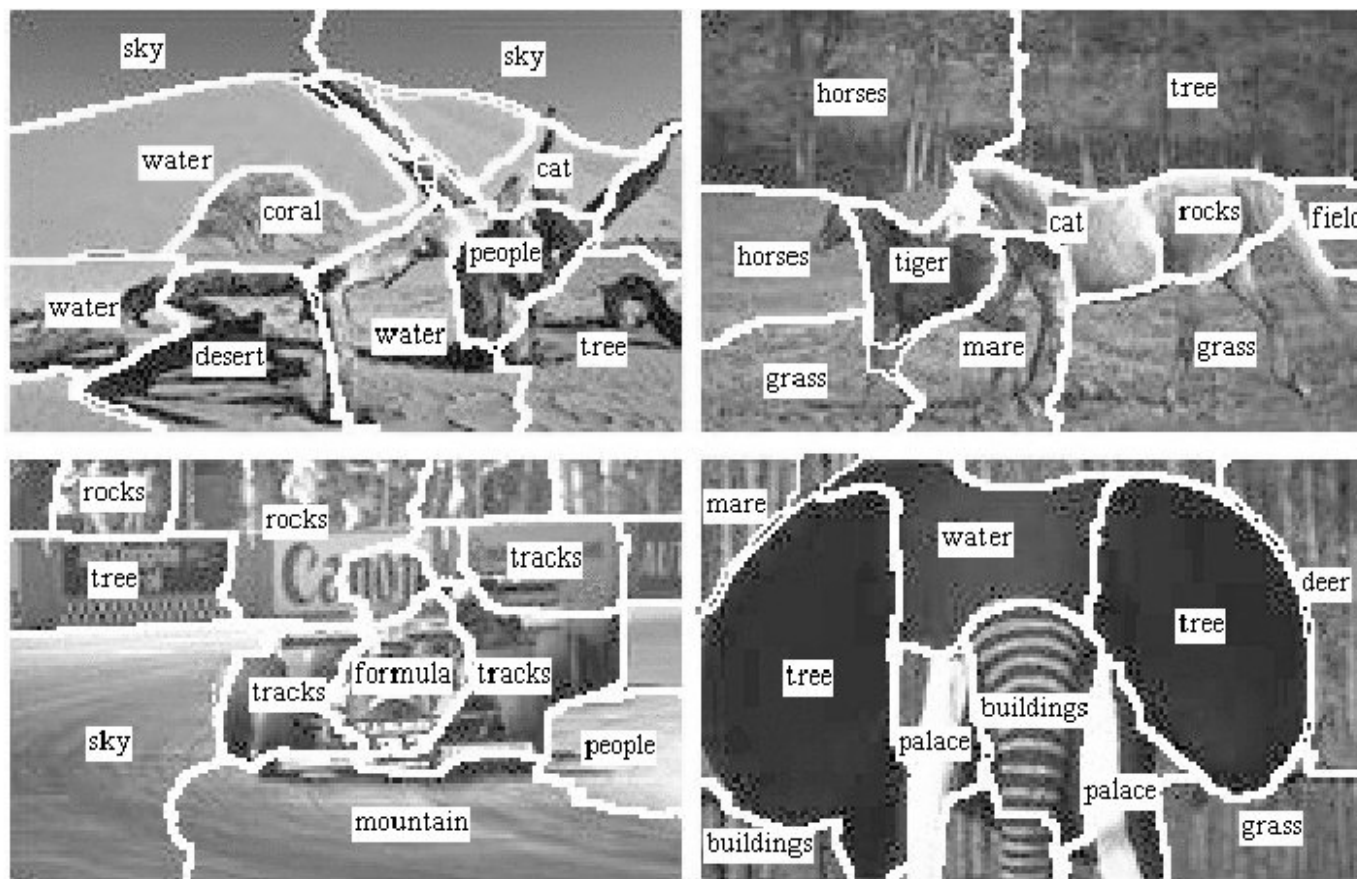
Matching words and Pictures - leveraging context

ⁿ **Approach 2** – paired word and region emission at nodes

$$p(D|d) = \sum_c p(c) \prod_{(w,b) \in D} \left[\sum_l p((w,b)|l,c) p(l|d) \right]$$

$$p(w \Leftrightarrow b) \approx \sum_c p(c) \sum_l p((w,b)|l,c) p(l|d).$$

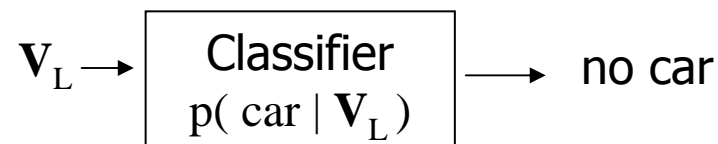
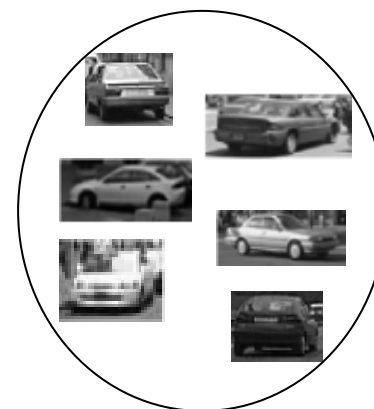
Some more results



Barnard, K., Duygulu, P., Forsyth, D., de Freitas, N., Blei, D. M., and Jordan, M. I. 2003. Matching words and pictures. *J. Mach. Learn. Res.* 3 (March 2003), 1107-1135.

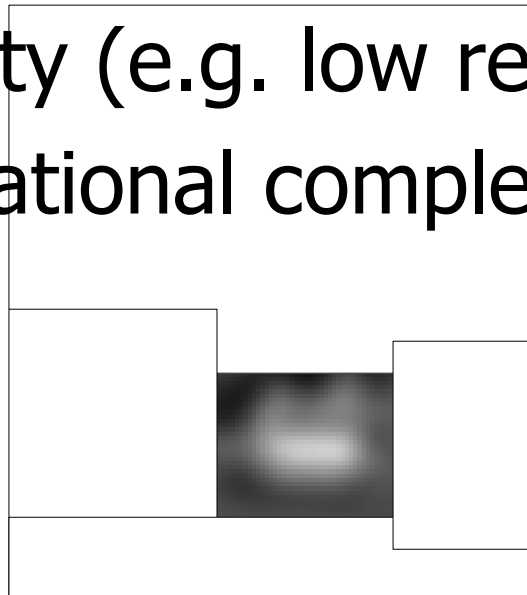
Traditional Approach to Object Detection

- n Uses only local object properties
- n Classify local image patches at each location and scale.



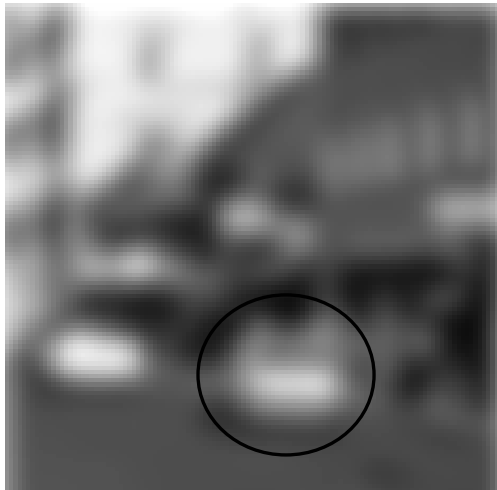
Problems of Local Methods

- n Ambiguity (e.g. low resolution)
- n Computational complexity



What's this?

Resolving Ambiguity



With contextual information, you can tell what it is.

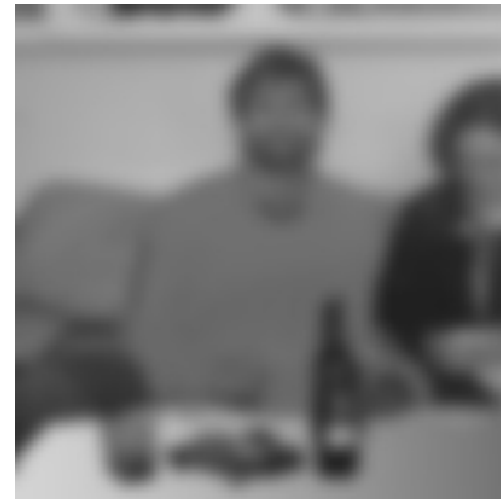
Resolving Ambiguity



pedestrian



car

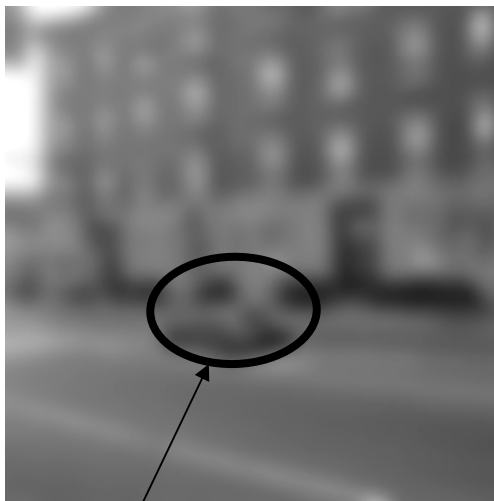


ash tray

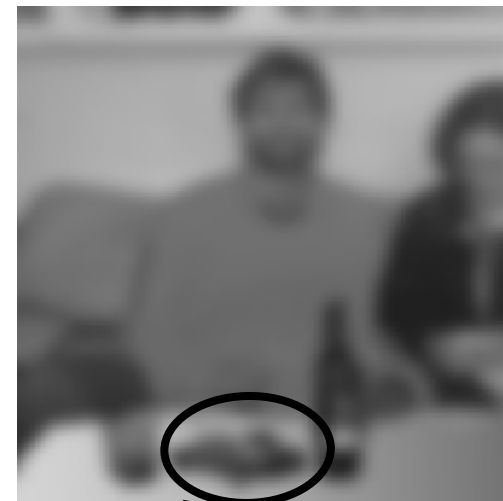
Resolving Ambiguity



pedestrian



car



ash tray

Identical local image features!



Context – Reduces Search Space

- n Context can provide a prior on what to look for, and where to look for it.
 - n Example
 - n Office → desk, chair, computers
 - n Outdoor → cars are usually on the road

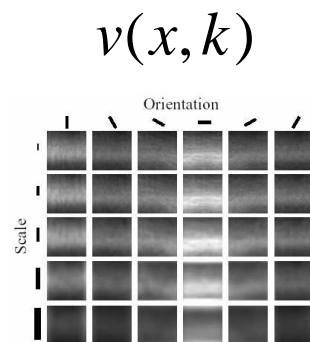


Context Features: Holistic Representation

- n Spatial layout of spectral components
 - n Responses of filters with different scales and orientations to image
 - n Filter examples
 - n Steerable Pyramid (Torralba et al., 2003)
 - n Weighted Fourier Transform (Torralba & Sinha 2001)
 - n Gabor filter banks (Torralba 2003)
 - n PCA dimensionality reduction

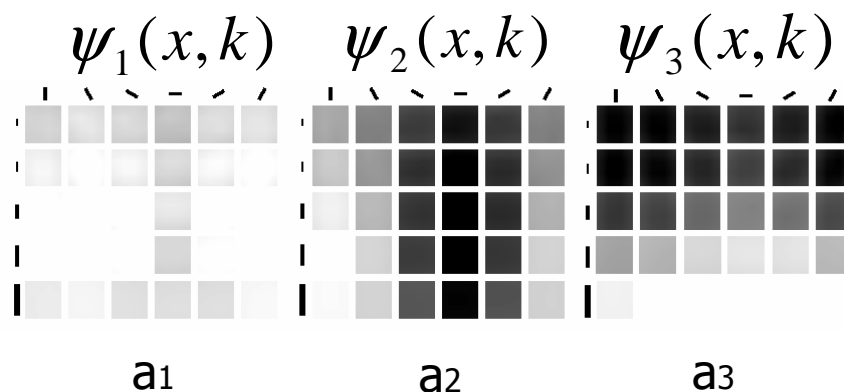
Context Features

After filtering...



$$v(x, k) \cong \sum_{n=1}^D a_n \cdot \psi_n(x, k)$$

PCA →



Context features:

$$v^c = [a_{1 \sim N}]$$



Statistical Context Priming

- n Context can be a rich source of information about object's identity, location and scale.
- n Statistical framework for modeling the relationship between context and object properties.

Statistical Context Priming

We are interested in $P(\vec{p}, \sigma, \vec{x}, o_n | \vec{v})$

Local evidence $P_l(\vec{p}, \sigma, \vec{x}, o_n | \vec{v}_L)$

Context priming $P_c(\vec{p}, \sigma, \vec{x}, o_n | \vec{v}_C)$

$$P_c(\vec{p}, \sigma, \vec{x}, o_n | \vec{v}_C) = P_p(\vec{p} | \sigma, \vec{x}, o_n, \vec{v}_C) P_s(\sigma | \vec{x}, o_n, \vec{v}_C) P_f(\vec{x} | o_n, \vec{v}_C) P_o(o_n | \vec{v}_C)$$

Pose and
Shape Priming

Scale
Selection

Focus of
Attention

Object
Priming

Object Priming

- n Annotated database for learning
- n Learning by MOG

$$P(o / \mathbf{v}_c) = \frac{P(\mathbf{v}_c / o)P(o)}{P(\mathbf{v}_c)}$$

$$P(\mathbf{v}_c) = P(\mathbf{v}_c / o)P(o) + P(\mathbf{v}_c / \neg o)P(\neg o)$$

$$P(\mathbf{v}_c / o) = \sum_{i=1}^M b_i \cdot G(\mathbf{v}_c; \mu_i, \Sigma_i) \quad P(\mathbf{v}_c / \neg o) = \sum_{i=1}^M b_i \cdot G(\mathbf{v}_c; \mu_i, \Sigma_i)$$

Model parameters $(b_i, \mu_i, \Sigma_i)_{i=1,M}$ are learned from **EM algorithm**.

Object Priming Result



Figure 7. Random selection of images from the test set showing the results of object priming for four superordinate object categories (o_1 = people, o_2 = furniture, o_3 = vehicles and o_4 = trees). The bars at the right-hand of each picture represent the probability $P(o | \mathbf{v}_C)$.

Object Priming Result

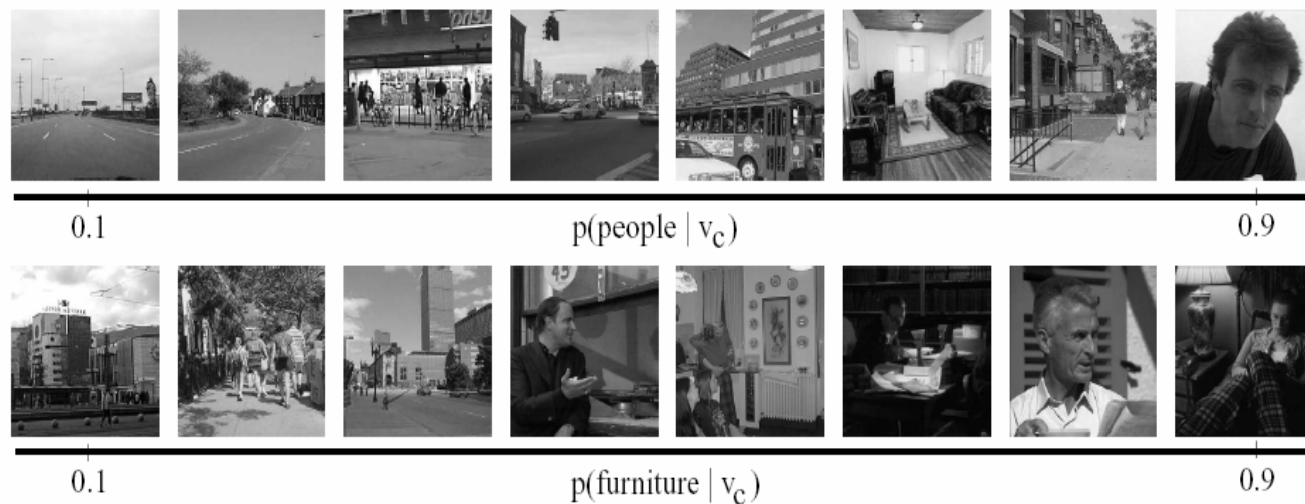


Fig. 8. Random selection of images from the test set organized with respect to the probability $P(o | v_C)$ for $o = \text{people}$ and furniture.

Context and Task-driven FOA

n Model: Gaussian Mixture Model

$$P(\mathbf{x} | o, \mathbf{v}_C) = \frac{\sum_{i=1}^M b_i G(\mathbf{x}; \mathbf{x}_i, \mathbf{X}_i) G(\mathbf{v}_C; \mathbf{v}_i, \mathbf{V}_i)}{\sum_{i=1}^M b_i G(\mathbf{v}_C; \mathbf{v}_i, \mathbf{V}_i)} \quad \mathbf{x}_i = \mathbf{a}_i + \mathbf{A}_i(\mathbf{v}_C - \mathbf{v}_i)$$

$$\begin{aligned} (\bar{x}, \bar{y}) &= \int \mathbf{x} P(\mathbf{x} | o, \mathbf{v}_C) d\mathbf{x} \\ &= \frac{\sum_{i=1}^M b_i \mathbf{x}_i G(\mathbf{v}_C; \mathbf{v}_i, \mathbf{V}_i)}{\sum_{i=1}^M b_i G(\mathbf{v}_C; \mathbf{v}_i, \mathbf{V}_i)} \quad \sigma_r^2 = \int r^2 P(\mathbf{x} | o, \mathbf{v}_C) d\mathbf{x} \end{aligned}$$

n Model learning by EM

$$(b_i, a_i, \mathbf{A}_i, \mathbf{X}_i, \mathbf{v}_i, \mathbf{V}_i)_{i=1,M}$$

Context and Task-driven FOA

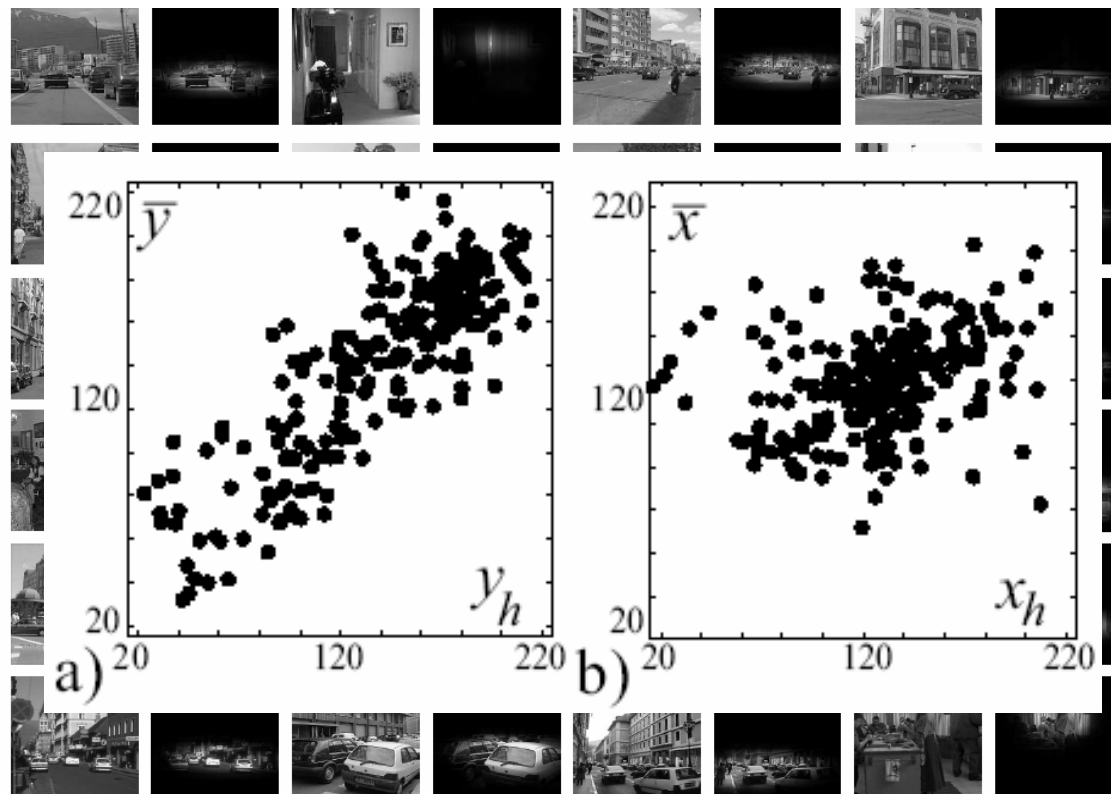


Fig. 14. Focus of attention based on global context configuration. Each pair shows the original image and the image multiplied by the function $P(x | v_C, o = heads)$ to illustrate the primed regions.

A. Torralba, "Contextual priming for object detection," IJCV 2003.

Context-driven Scale Selection

n Model (GMM)

$$P(\sigma | o, \mathbf{v}_C) = \frac{\sum_{i=1}^M b_i G(\sigma; \sigma_i, \mathbf{S}_i) G(\mathbf{v}_C; \mathbf{v}_i, \mathbf{V}_i)}{\sum_{i=1}^M b_i G(\mathbf{v}_C; \mathbf{v}_i, \mathbf{V}_i)} \quad \sigma_i = a_i + \mathbf{A}_i(\mathbf{v}_C - \mathbf{v}_i)$$

$$\bar{\sigma} = \int \sigma P(\sigma | o, \mathbf{v}_C) d\sigma = \frac{\sum_{i=1}^M \sigma_i b_i G(\mathbf{v}_C; \mathbf{v}_i, \mathbf{V}_i)}{\sum_{i=1}^M b_i G(\mathbf{v}_C; \mathbf{v}_i, \mathbf{V}_i)}$$

$$\sigma_h^2 = \int (\sigma - \bar{\sigma})^2 P(\sigma | o, \mathbf{v}_C) d\sigma$$

n Model Learning by EM

$$(b_i, a_i, \mathbf{A}_i, \mathbf{S}_i, \mathbf{v}_i, \mathbf{V}_i)_{i=1,M}$$

Context-Driven Scale Selection

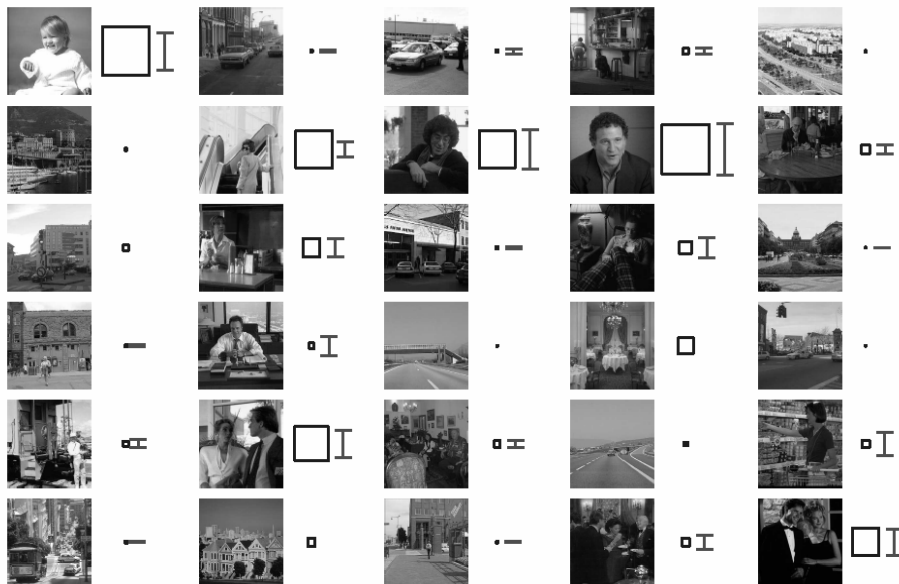


Fig. 19. Results for scale selection given global context information for rand

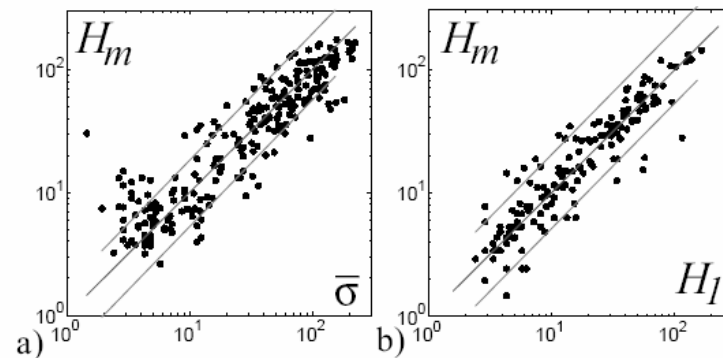


Fig. 22. Each row shows a set of 8 pictures sorted according to the predicted size of human heads (top) and cars (bottom). A. Torralba, "Contextual priming for object detection," IJCV 2003.



Context-based Place and Object Recognition

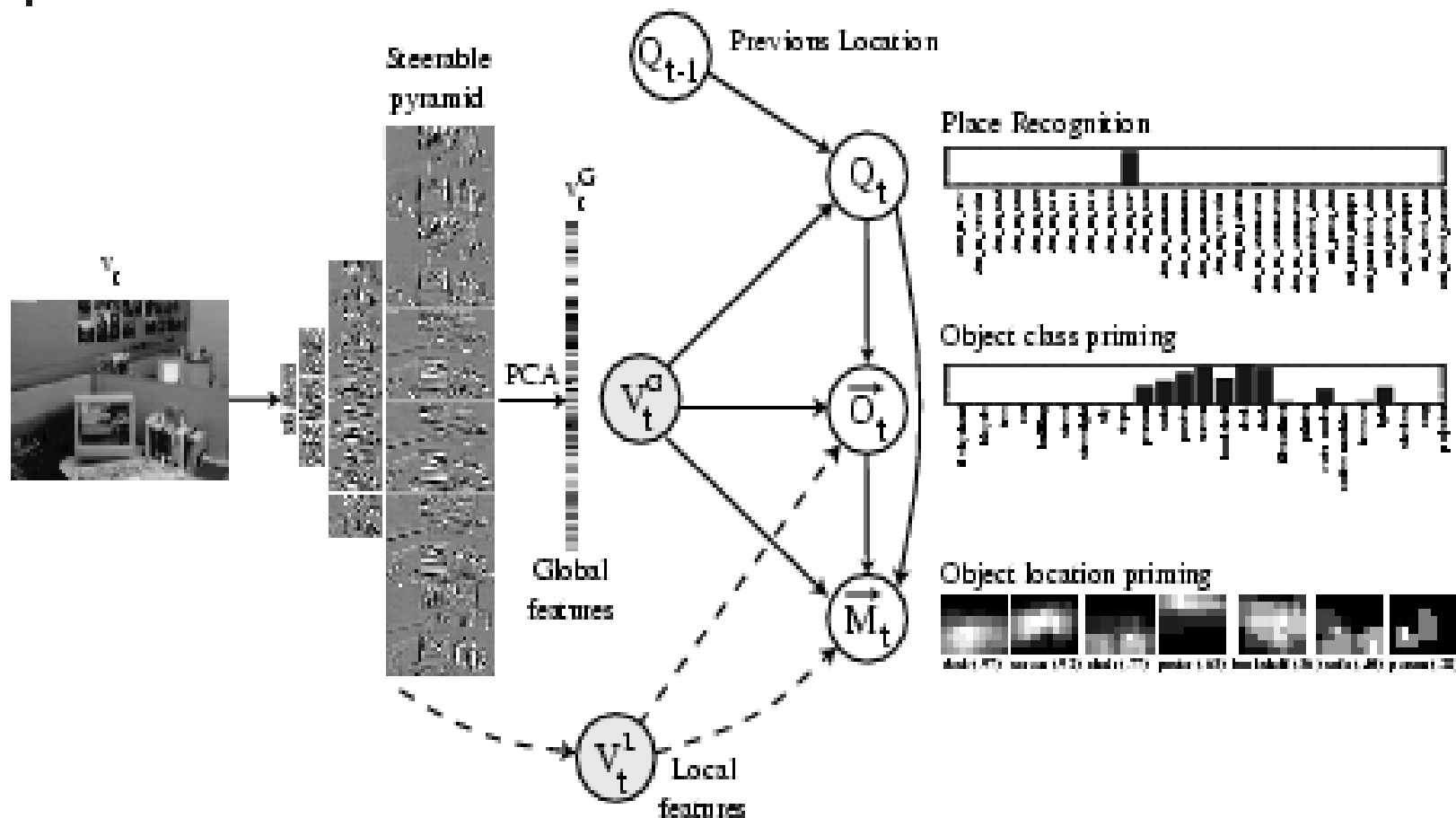
- n Testbed
 - n Wearable device: helmet-mounted mobile system
- n Place and Scene Recognition
 - n Global features + PCA
 - n Wearable test-bed
 - n Model of place recognition – HMM
- n Object Recognition
 - n Individual object recognition using contextual information



Place, Category Recognition

- n Place: Instance of a scene category
 - n E.g. office 610, main street
- n Scene category: type of place
 - n E.g. office, street

System Framework



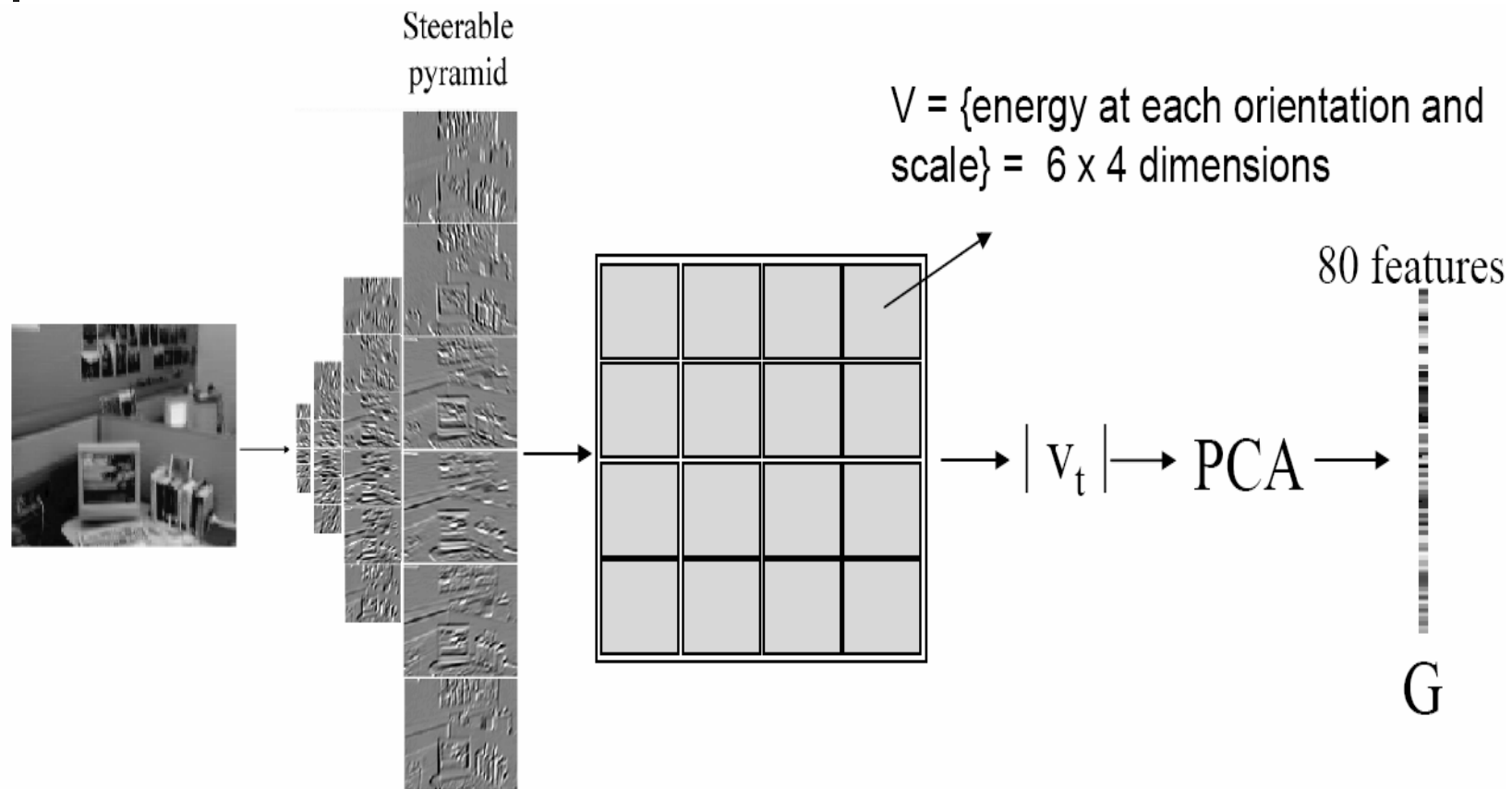
A. Torralba et al., "Context-based vision system for place and object recognition," ICCV 2003.



Low-dimensional Scene Representation

- Compute image intensity (no color)
- Steerable pyramid (6 orientations, 4 scales)
- Compute magnitude of average filter responses
- Downsample to 4×4 384 dimensions
- PCA to 80 dimensions

Global Feature (Interpretation)



Visualizing the Global Features

Images



Gist: 80-dim
Representation



A. Torralba et al., "Context-based vision system for place and object recognition," ICCV 2003.

Models for Place Recognition

n Transition model: HMM

$$\begin{aligned} P(Q_t = q | v_{1:t}^G) &\propto p(v_t^G | Q_t = q) P(Q_t = q | v_{1:t-1}^G) \\ &= p(v_t^G | Q_t = q) \sum_{q'} A(q', q) P(Q_{t-1} = q' | v_{1:t-1}^G) \end{aligned}$$

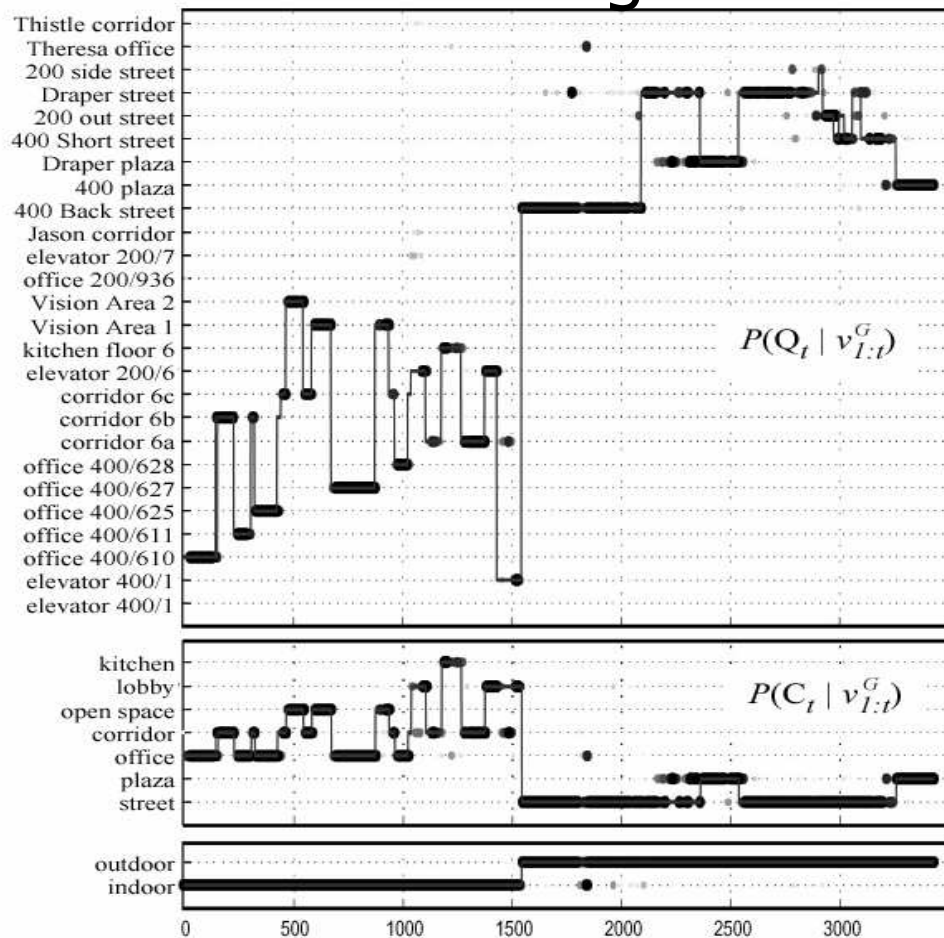
$$A(q', q) = P(Q_t = q | Q_{t-1} = q')$$

n Observation model $p(v_t^G | Q_t)$

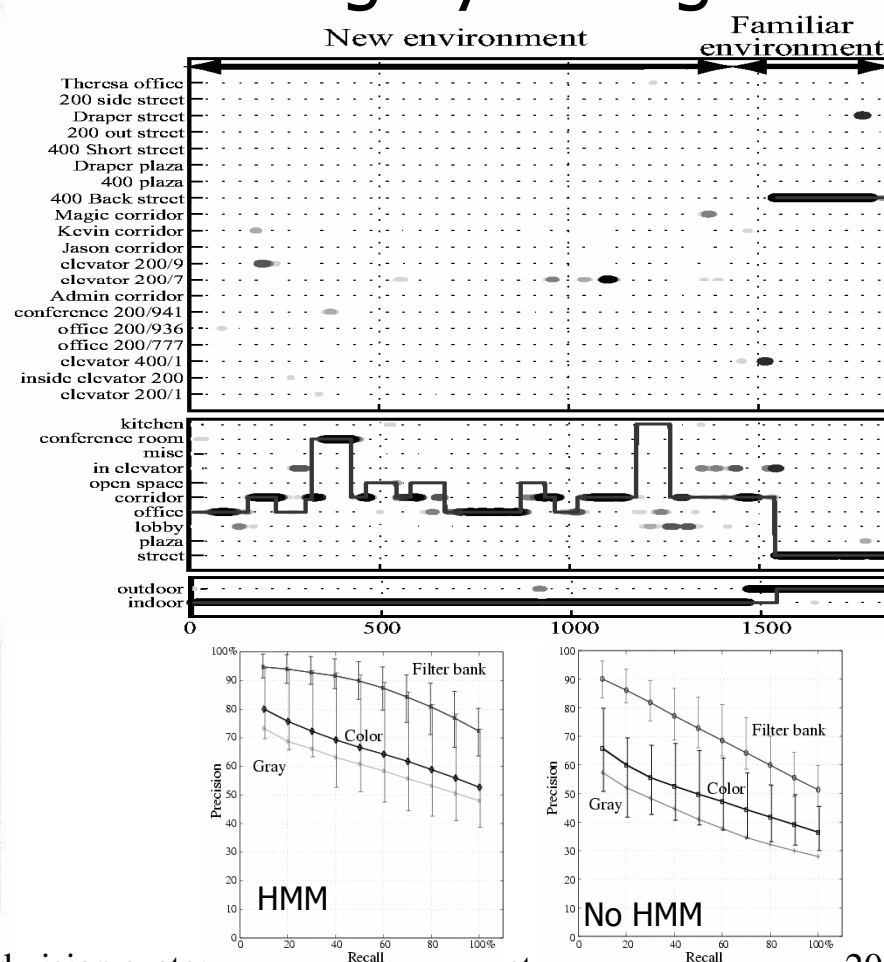
n Mixture of Gaussian (100views/Places)

Recognition Performance

Place Recognition



Category Recognition



A. Torralba et al., "Context-based vision system for place and object recognition," ICCV 2003.

Contextual Priming for Object Detection

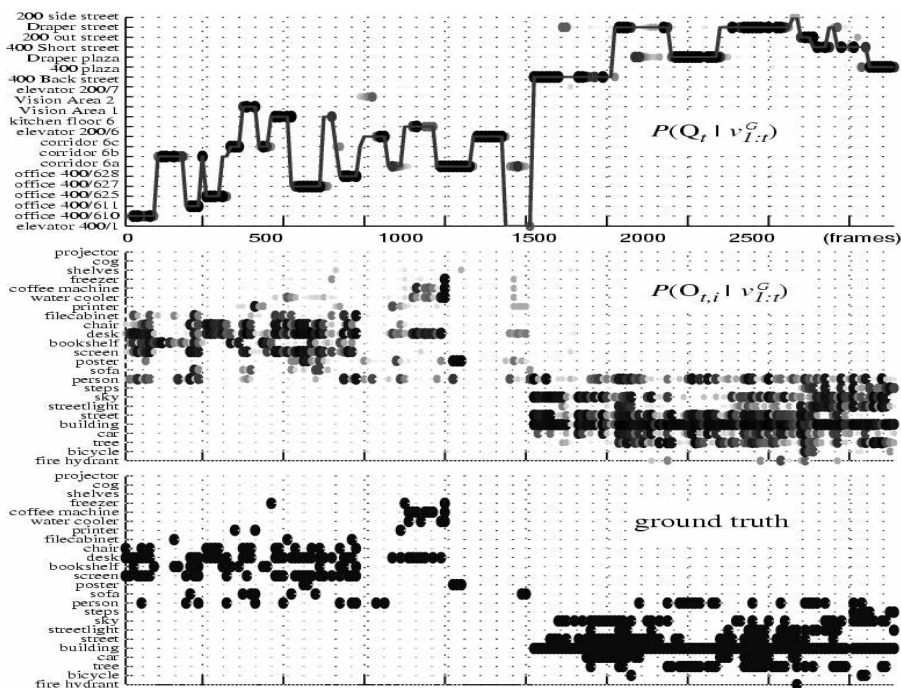
Object properties

$$\vec{O}_t = (O_{t,1}, \dots, O_{t,N_o})$$

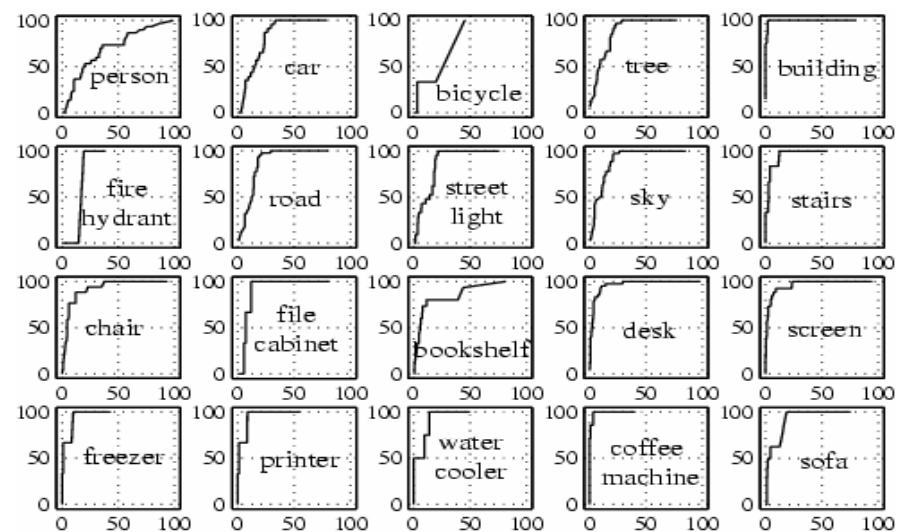
Object priming using context ONLY

$$P(O_{t,i}|v_t^G, Q_t = q) = \frac{p(v_t^G|O_{t,i}, j)P(O_{t,i}|q)}{p(v_t^G|O_{t,i}, q)P(O_{t,i}|q) + p(v_t^G|\bar{O}_{t,i}, q)P(\bar{O}_{t,i}|q)}$$

Predicting Object Presence



Object presence prob.



Perf. by ROC curve

Contextual Priors for Object Localization

n Goal: $P(X_{t,i}|v_t^G, O_{t,i} = 1)$

n Approach:

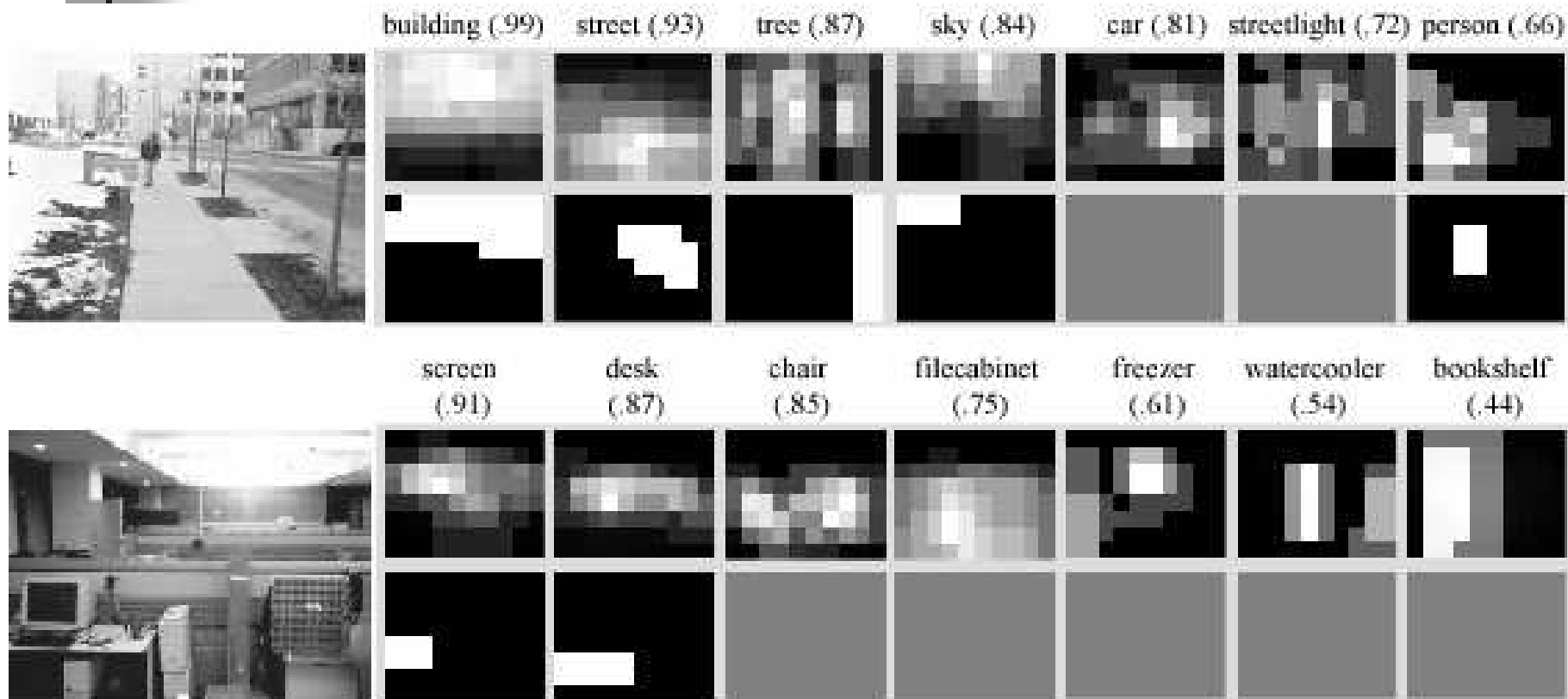
n Divide images into 8×10 binary masks

$$E[M_{t,i}|v_{1:t}^G] = \sum_{o \in \{0,1\}} \sum_q P(O_{t,i} = o, Q_t = q | v_{1:t}^G) \\ \times E[M_{t,i}|v_t^G, Q_t = q, O_{t,i} = o]$$

n Likelihood is weighted sum of prototypes

$$E[M_{t,i}|v_t^G] = \sum_q \sum_k w_{k,i,q} \times \mu_{k,i,q}^m$$

Results for Object Localization





Other Models of Context

- n Using forest to see the trees: graphical model, K. Murphy 2003.
- n Boosted random field, K. Murphy 2004.



Summary

- n We introduced basic concepts of context and its role in human visual perception.
- n We also discussed some recent probabilistic approaches to application of context in computer vision.
 - n Matching words and pictures
 - n Statistical context priming for object recognition
- n Context is important in object recognition, especially when visibility of object appearance is low. It reduces ambiguity, search space, and also make a complementary role to the local approach.



References

- n Antonio Torralba and Pawan Sinha, "Statistical context priming for object detection," ICCV 2001.
- n Antonio Torralba, "Contextual priming for object detection," IJCV 2003.
- n Antonio Torralba et al., "Context-based vision system for place and object recognition," ICCV 2003.
- n Kevin Murphy et al., "Using the forest to see the trees: a graphical model relating features, objects, and scenes," NIPS 2003.



Graphical Model Relating Features, Objects, and Scenes

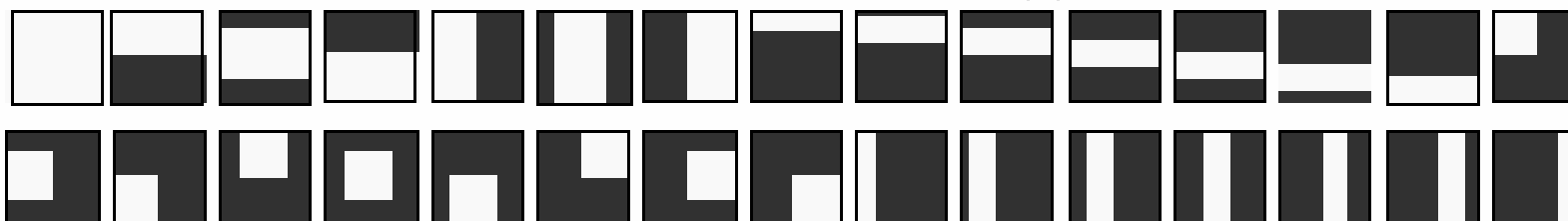
- n Combining local, bottom-up information with global, top-down information using a graphical model
 - n Boosting-based object detection
 - n Scene learning and classification
 - n Joint classification and detection

Filters and Spatial Templates

$$f_i(k) = \sum_x \tilde{w}_k(x) (|I(x) * g_k(x)|^{\gamma_k})_i \quad \gamma_k \in \{2, 4\}$$



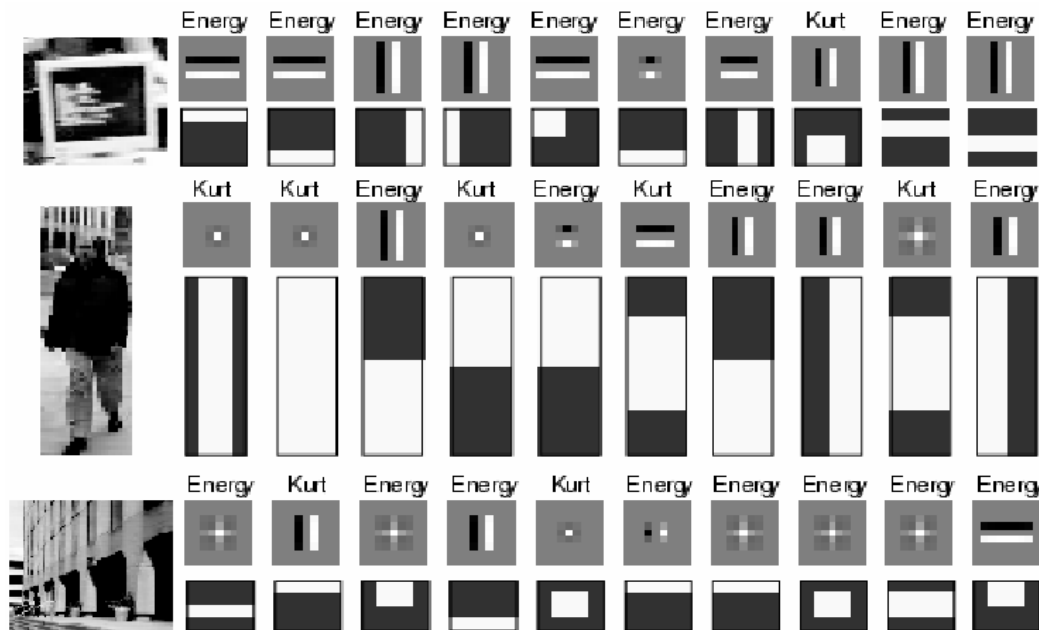
(a) Dictionary of 13 filters, $g(x)$.



(b) Dictionary of 30 spatial templates, $w(x)$.

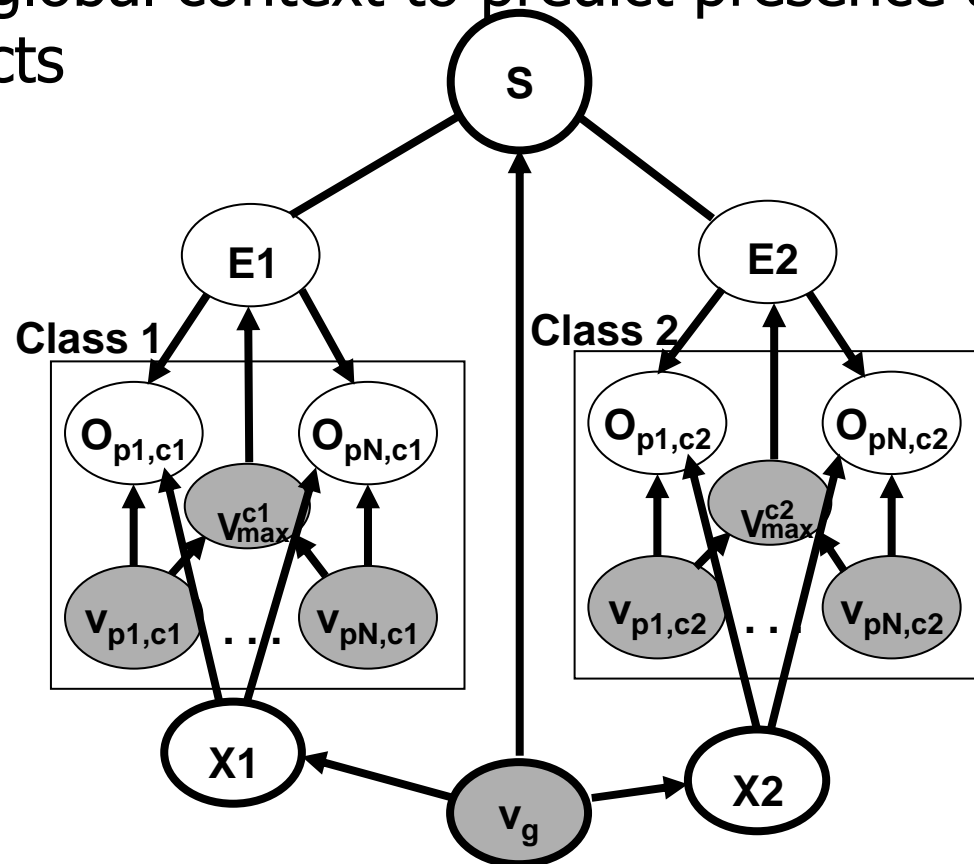
Examples of Learned Features

- Some features after 100 rounds of boosting



Tree-structured Graphical Model for Joint Classification and Detection

- n Use global context to predict presence and location of objects





Backup Slides

Steerable Pyramid

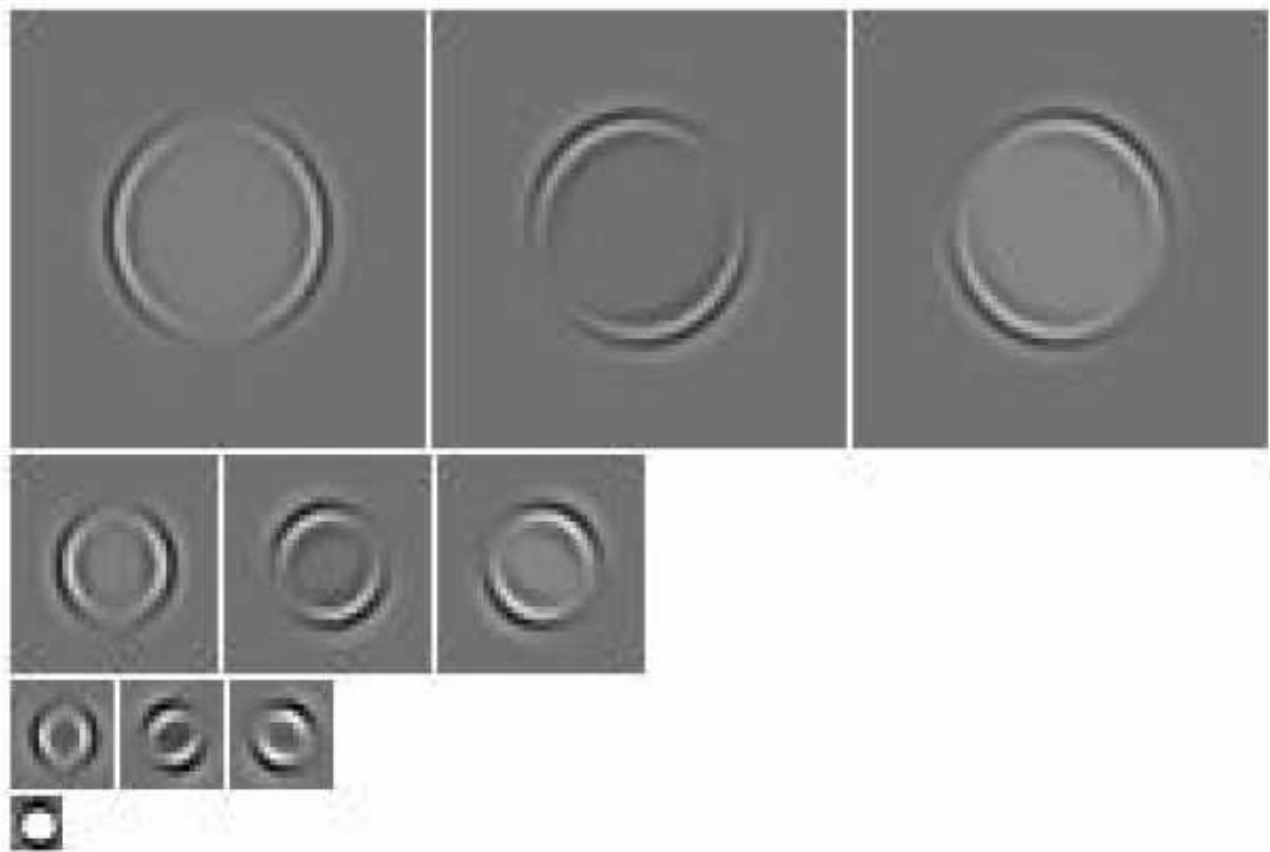


Image from Eero P Simoncelli & W.T. Freeman 1995

Weighted Fourier Transform

n WFT

$$I(x, y, f_x, f_y) = \sum_{x', y'=0}^{N-1} i(x', y') h_r(x' - x, y' - y) e^{-j 2\pi(f_x x' + f_y y')}$$

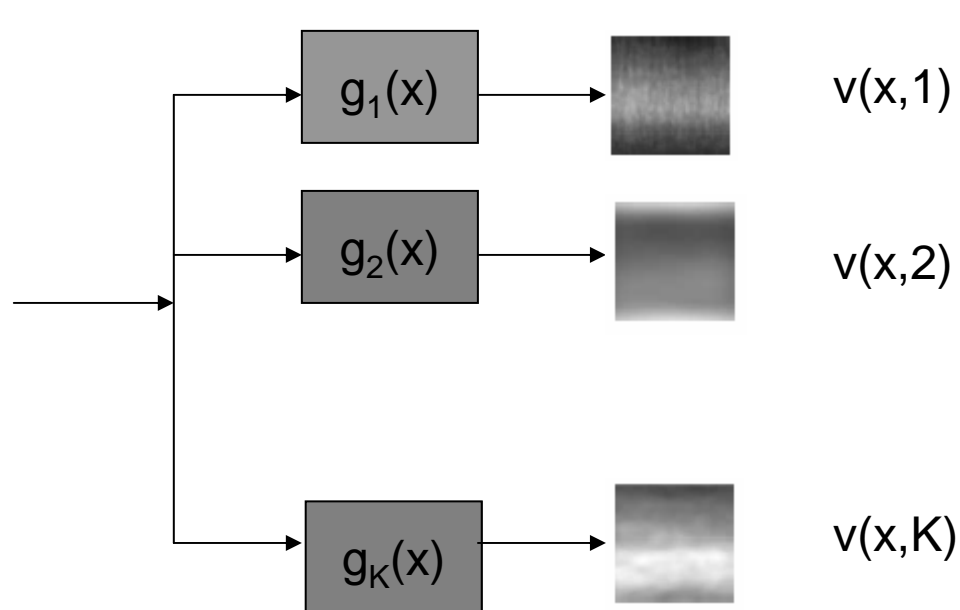
$$A(x, y, f_x, f_y) = \frac{|I(x, y, f_x, f_y)|}{I(x, y, 0, 0) \text{std}(x, y, f_x, f_y)} \quad \text{Normalization}$$

$$A(x, y, f_x, f_y) \simeq \sum_{n=1}^N a_n \psi_n(x, y, f_x, f_y) \quad \text{PCA with } N=60 \text{ PCs}$$

$$\boxed{\vec{v}_C = \{a_n\}_{n=1, N}} \longrightarrow \text{Context features}$$

Gabor Filter Bank

n Gabor filter bank



$$v(x, k) = \left| \sum_{x'} I(x') g_k(x - x') \right|$$

Magnitude of averaged filter output

$$g_k(x) = g_0 \cdot e^{-\frac{\|x\|^2}{\sigma_k^2}} \cdot e^{2\pi i \langle f_k, x \rangle}$$

Gabor (oriented band-pass) filters