

Natural Scene Recognition: From Humans to Computers

Li Fei-Fei

1. Computer Science Department
2. Psychology Department



A picture is worth a thousand words.

**--- Confucius
or *Printers' Ink* Ad (1921)**



blue

rugged

white and red

bright

textured structure

green

porous

grey

elongated shapes





- To **understand** human visual intelligence by via psychophysical and physiological experiments
- To **build** intelligent visual algorithms for machines and robots

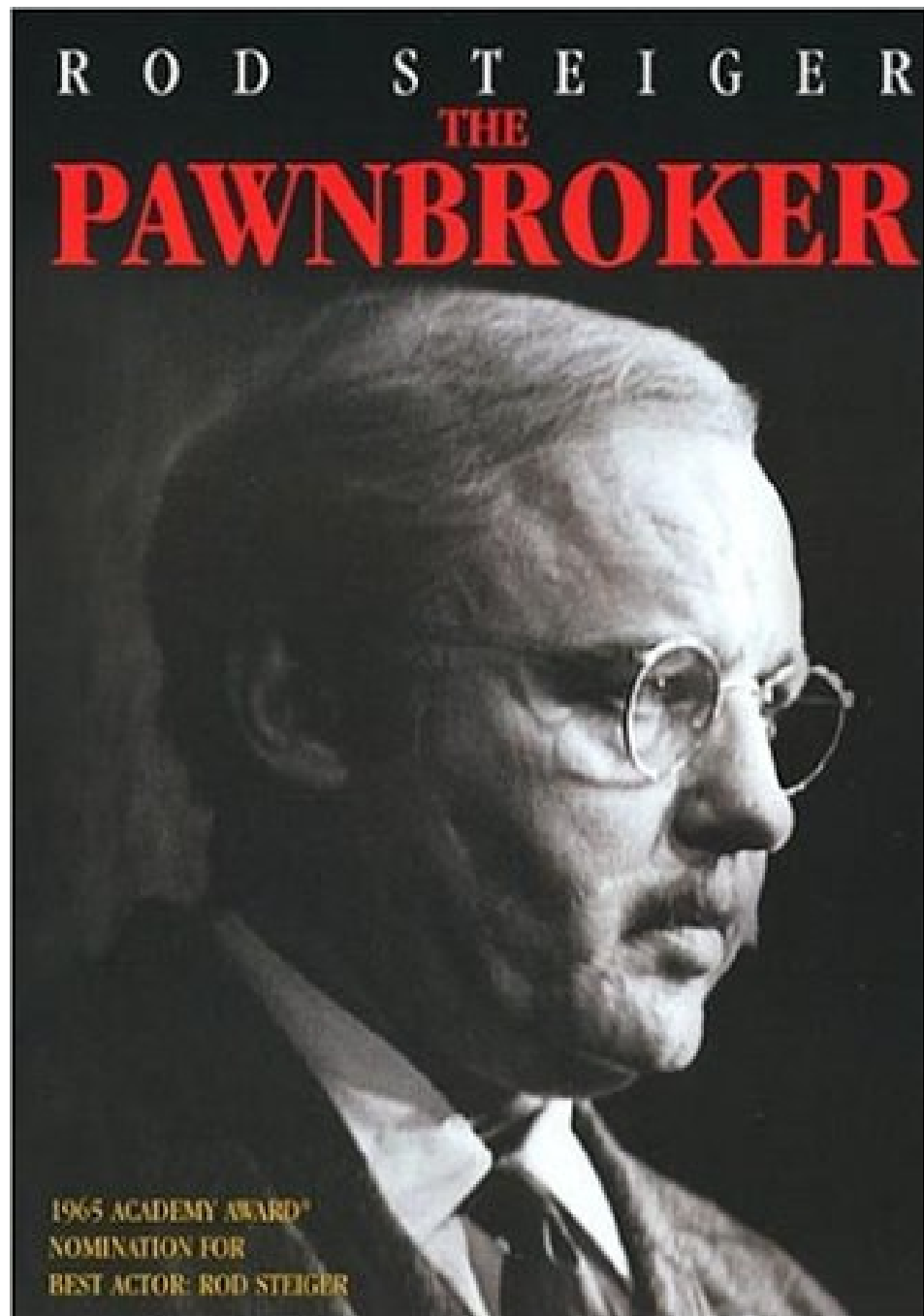
beach



living room

city

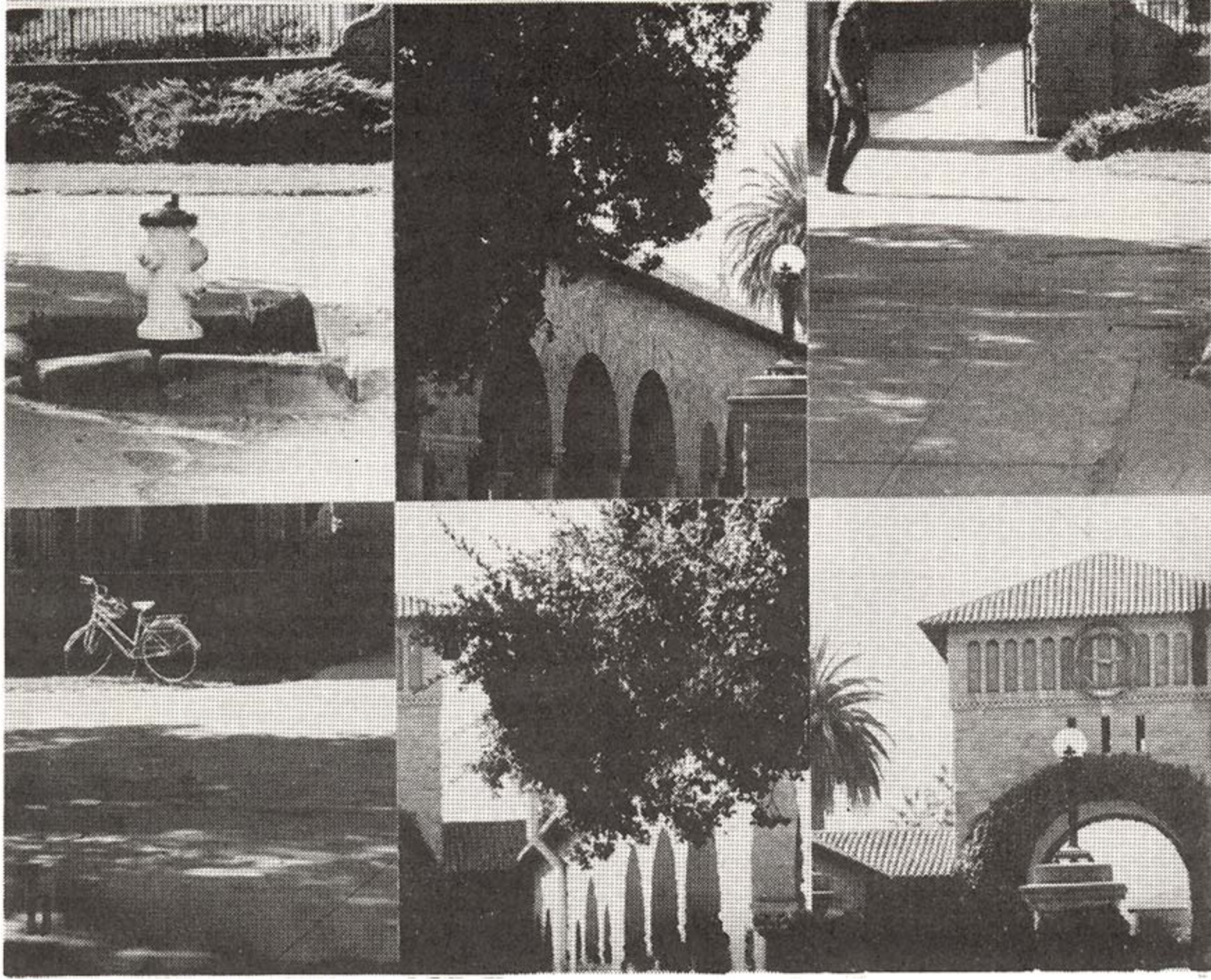




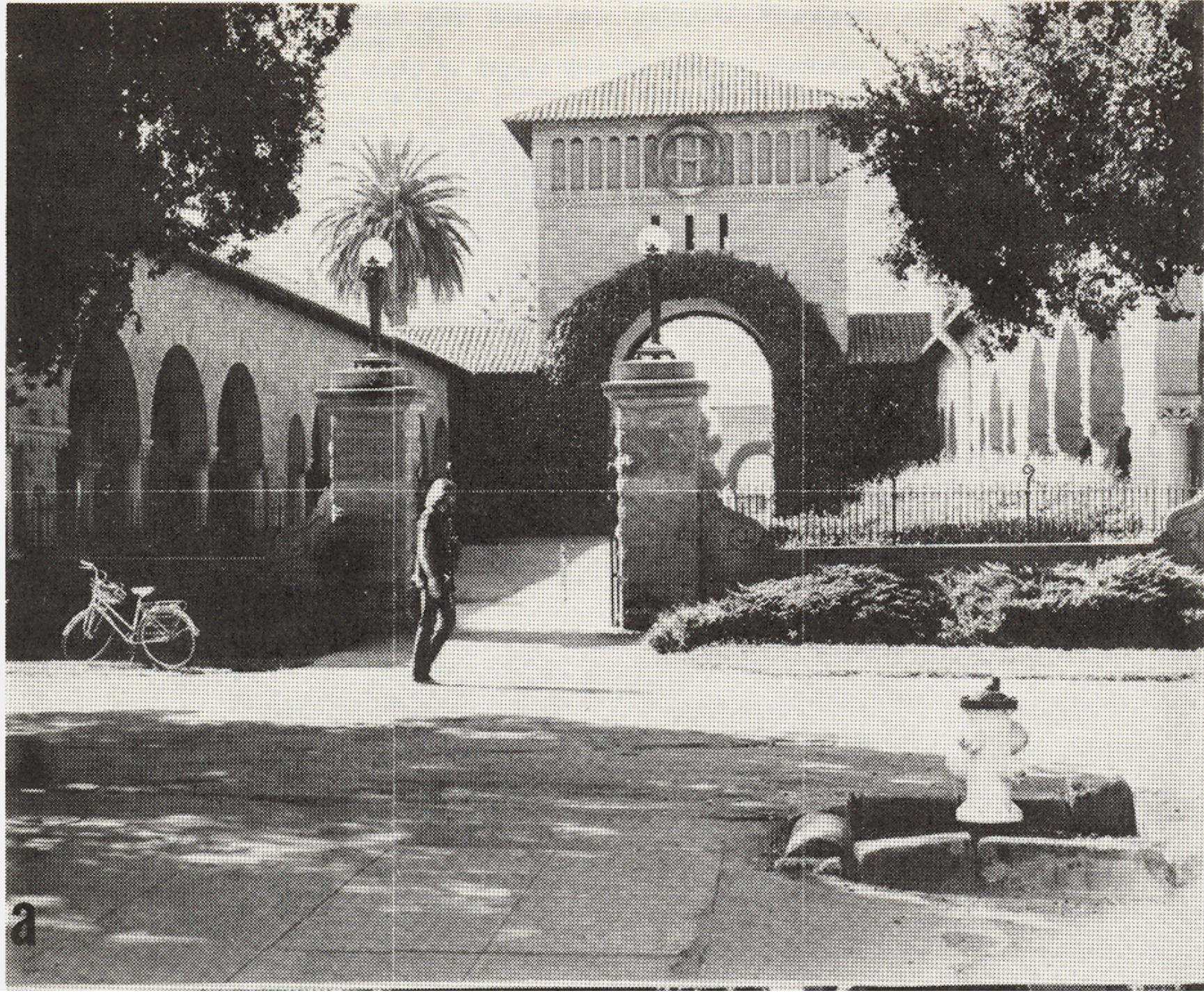
S. Lumet, 1965



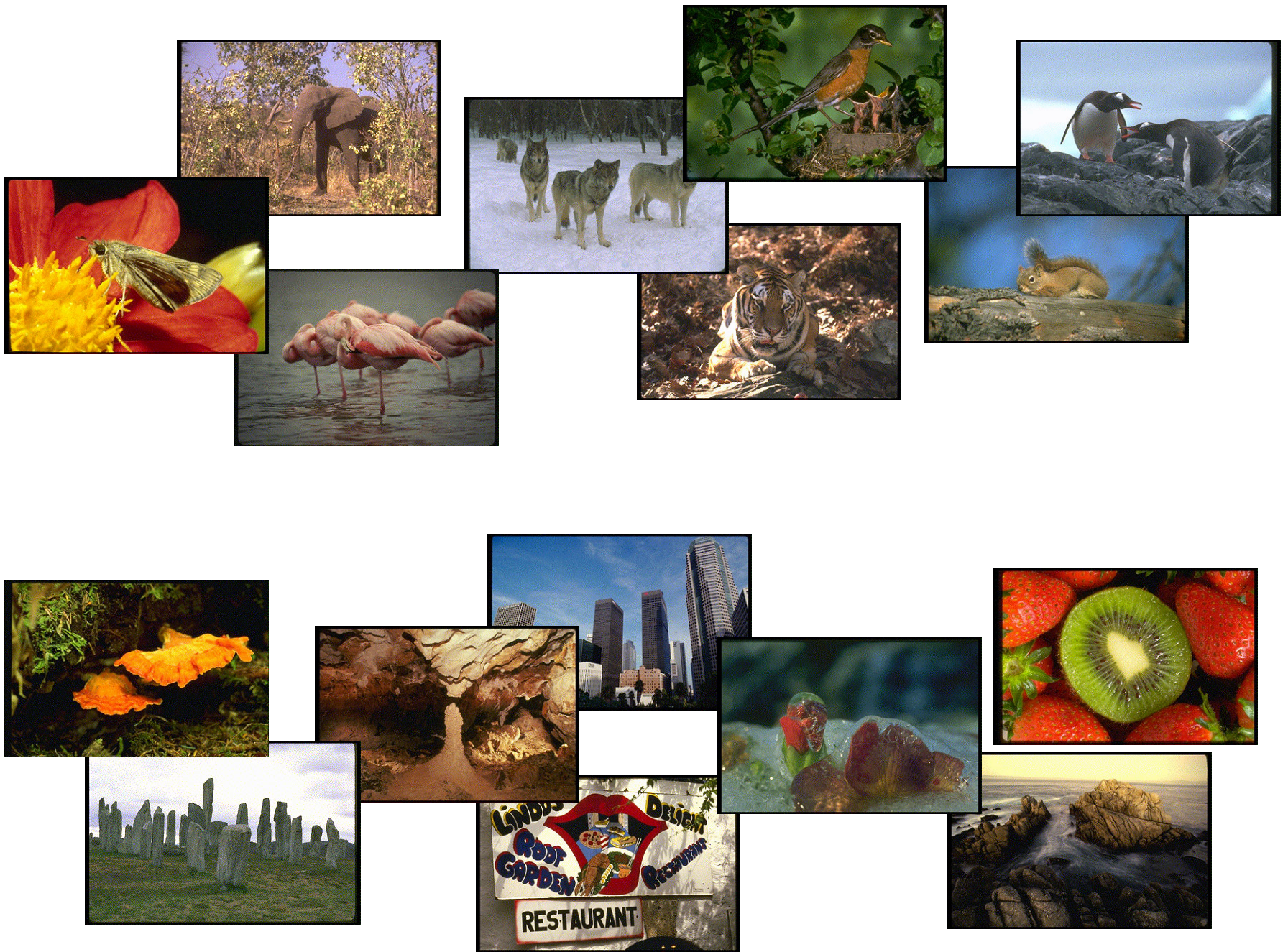
Potter, Biederman, etc. 1970s

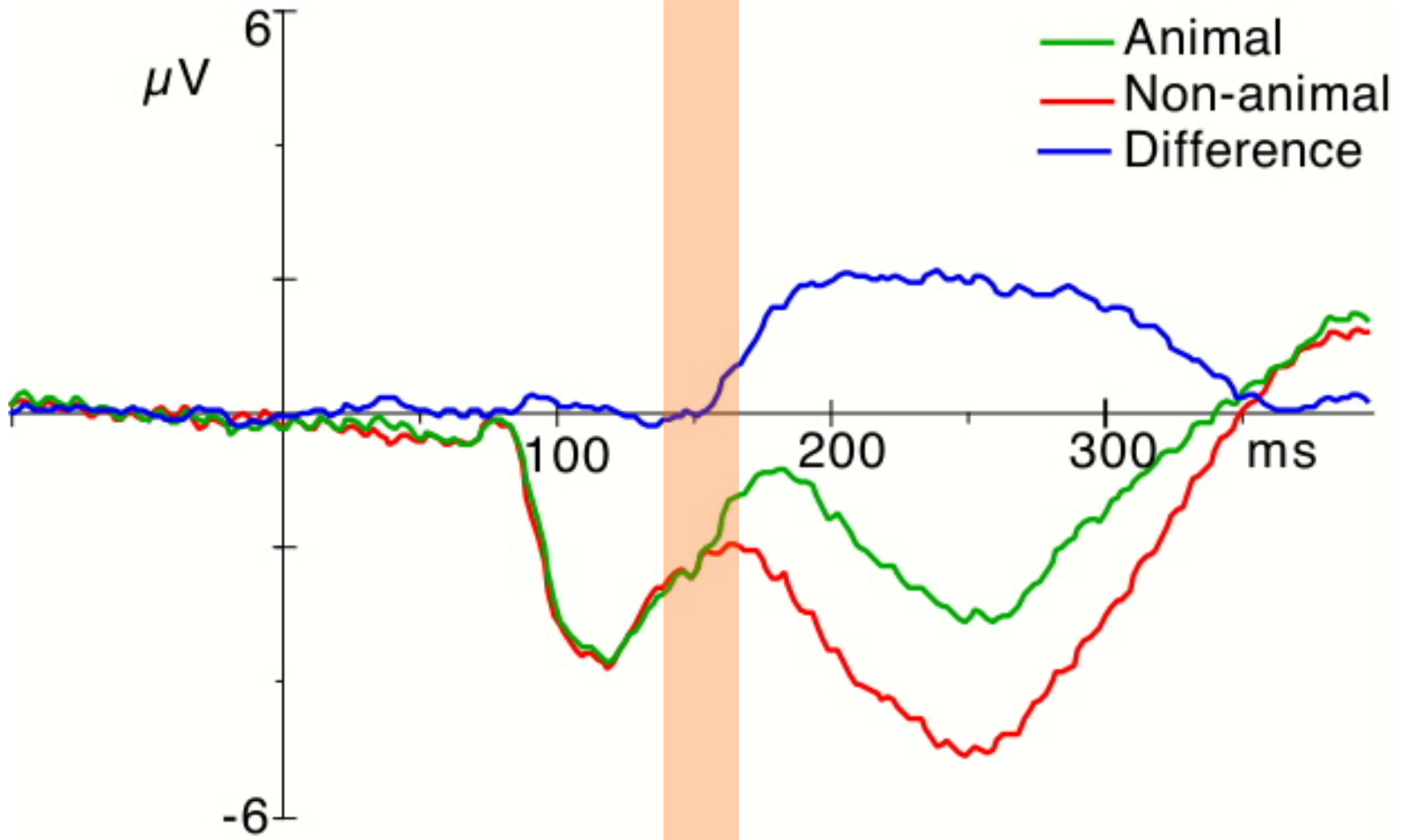


Biederman, *Science*, 1973

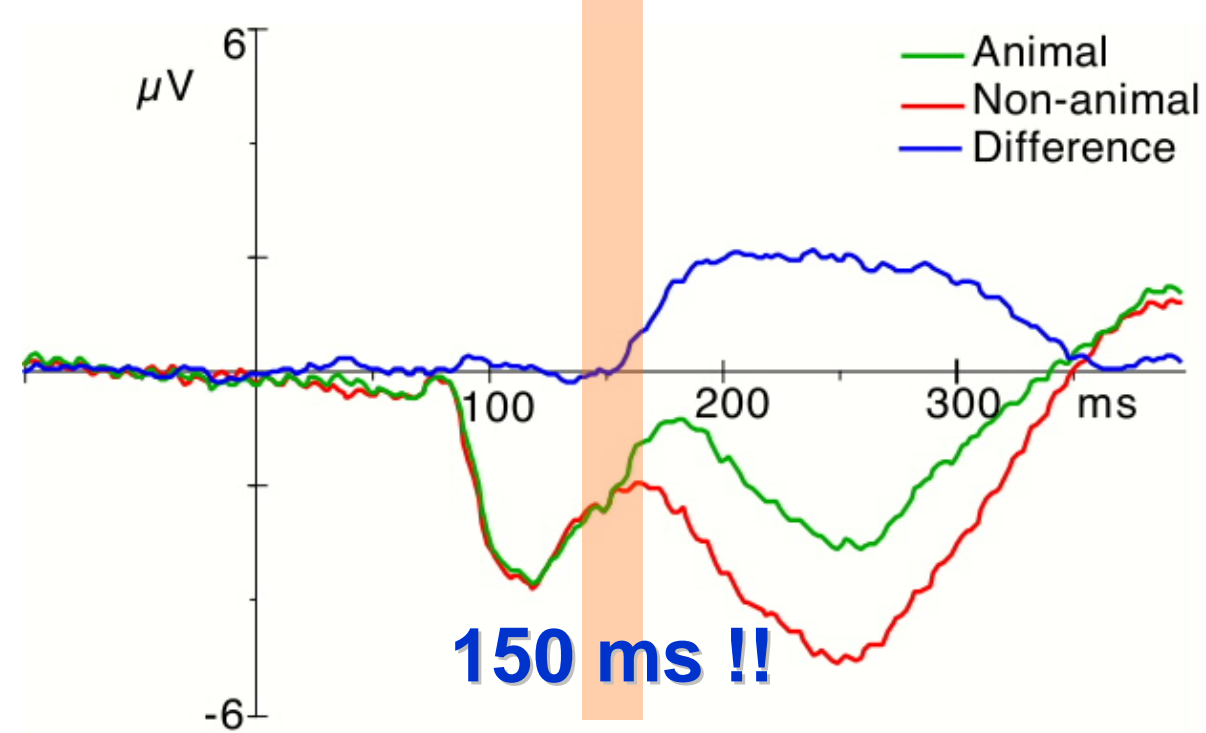


a

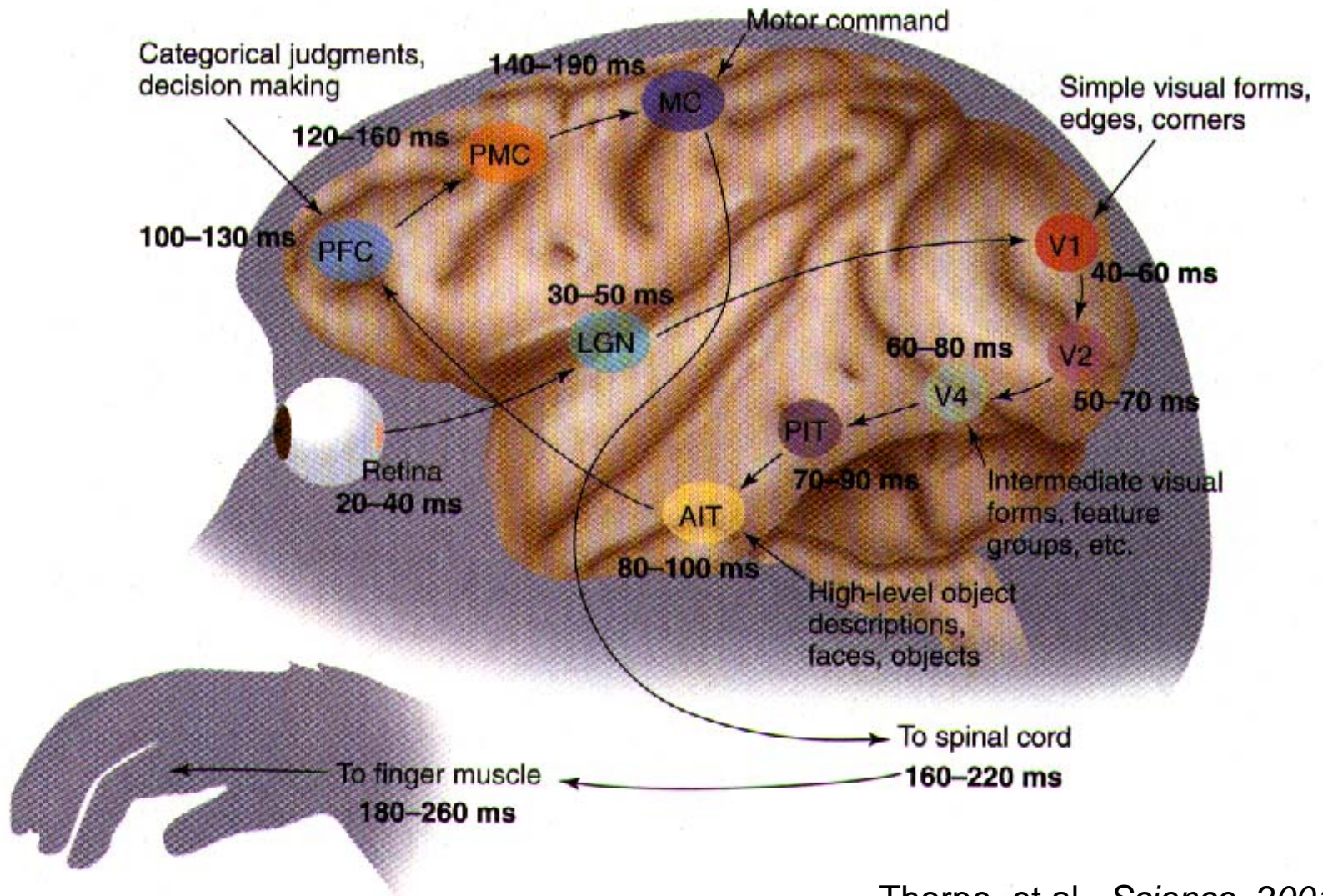




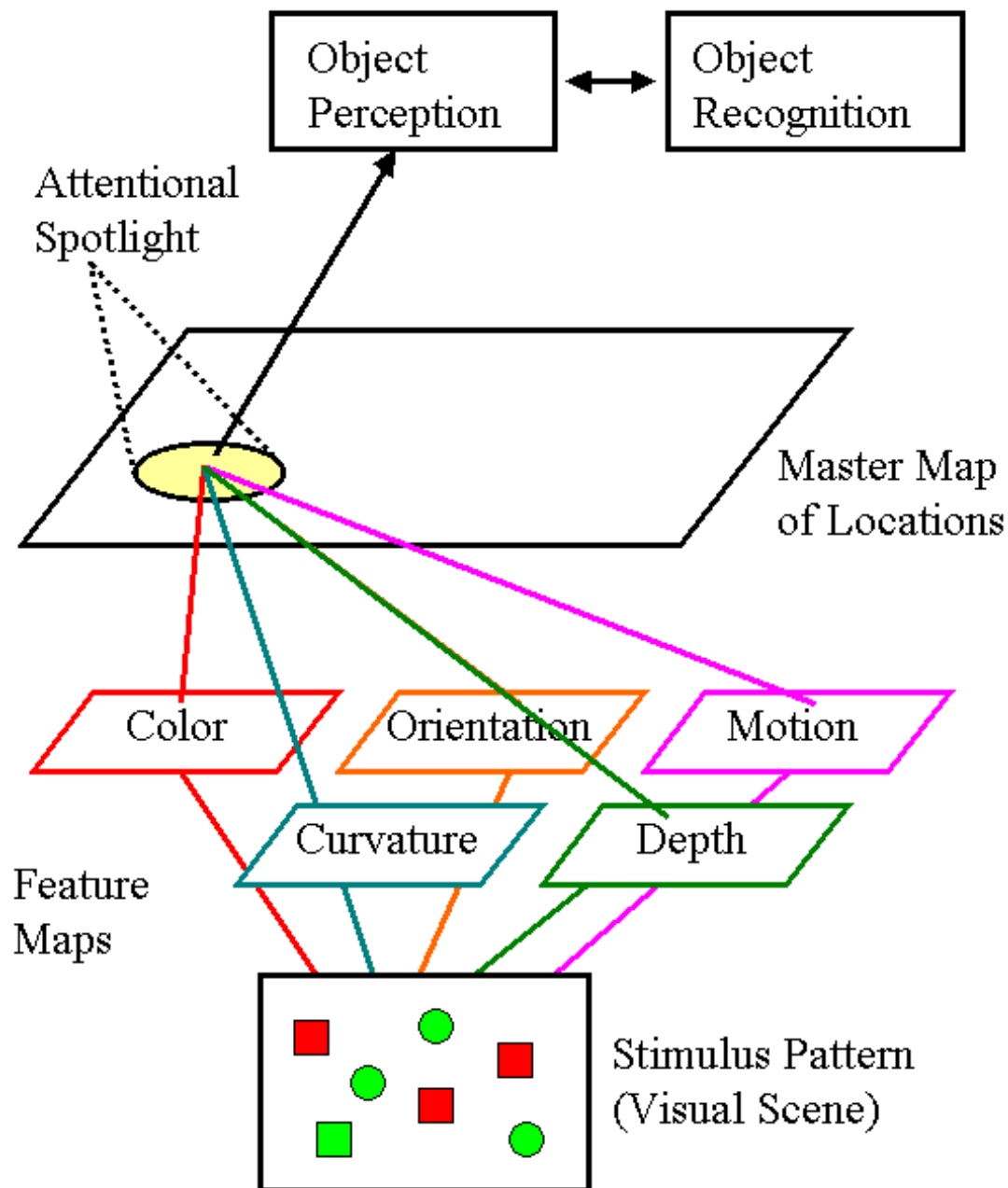
150 ms !!



A feed-forward mechanism?

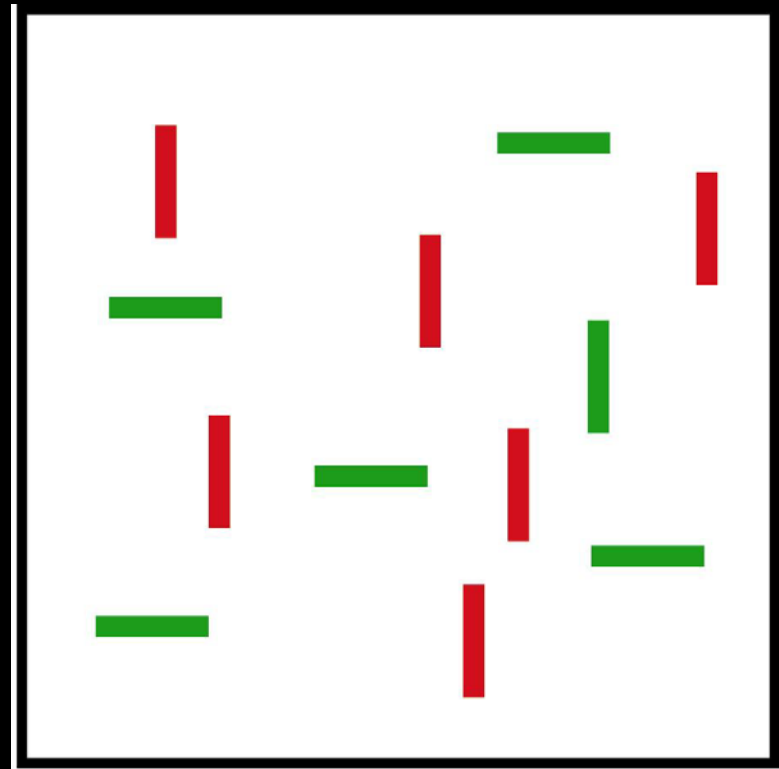
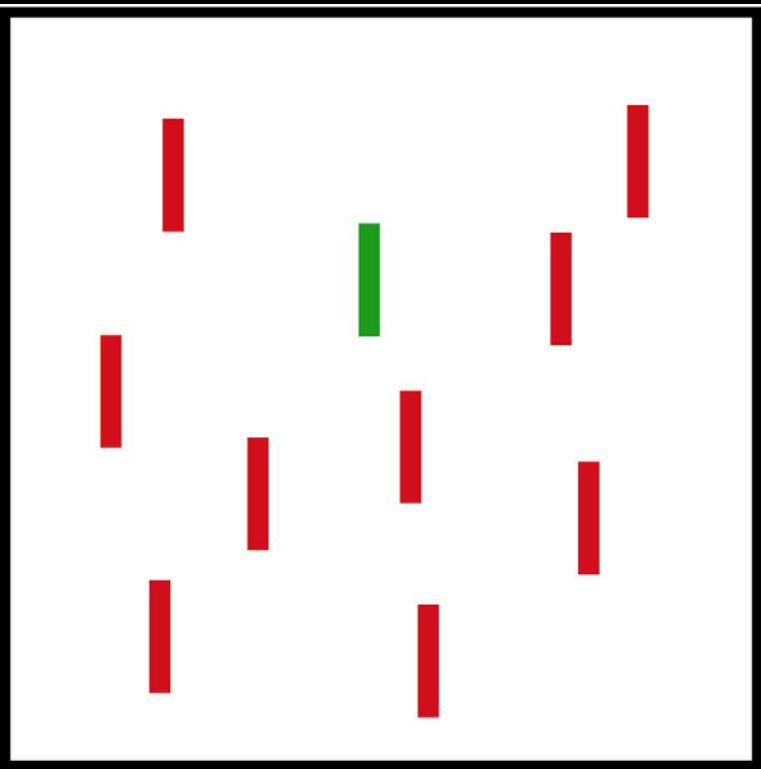


Feature Integration Theory

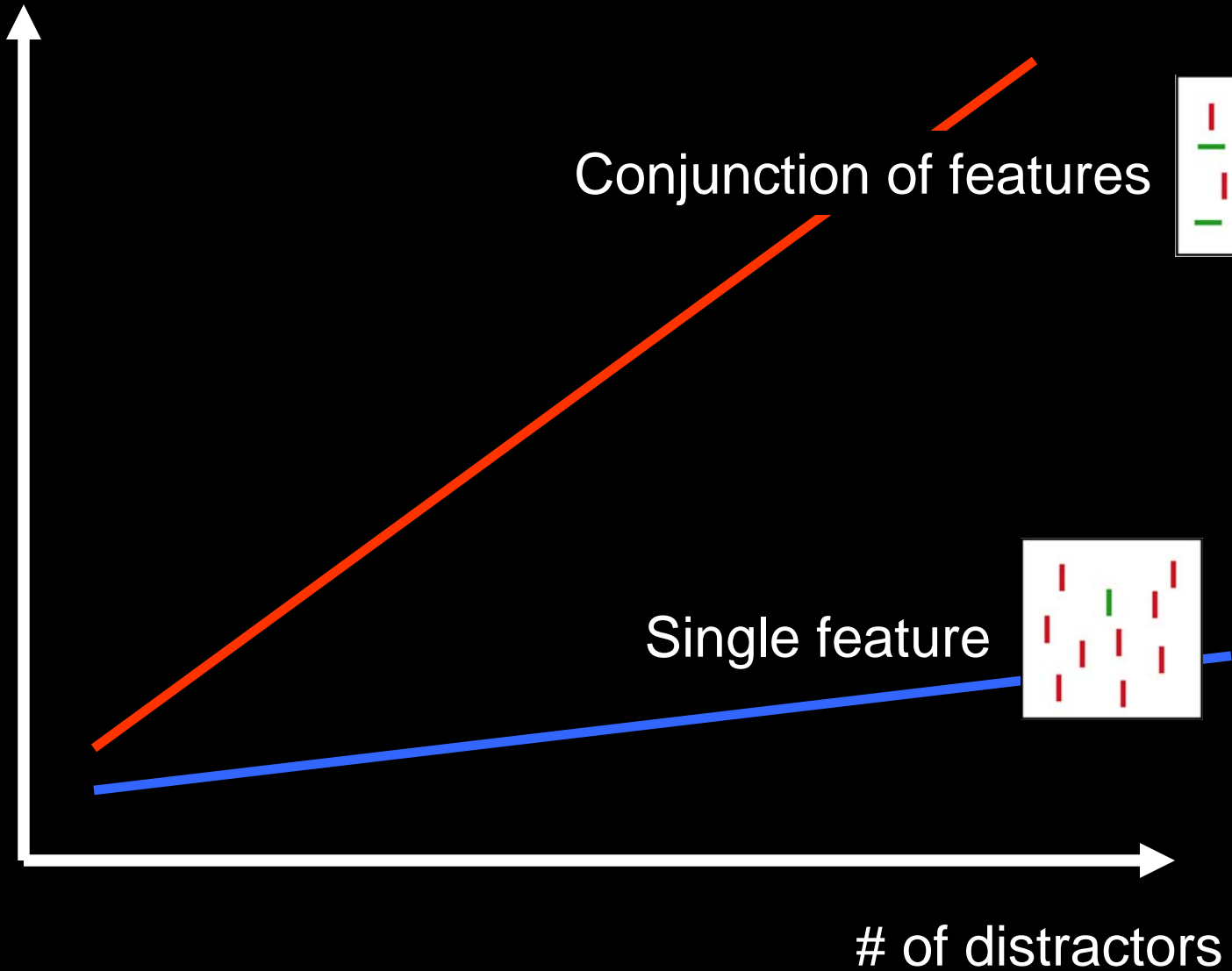


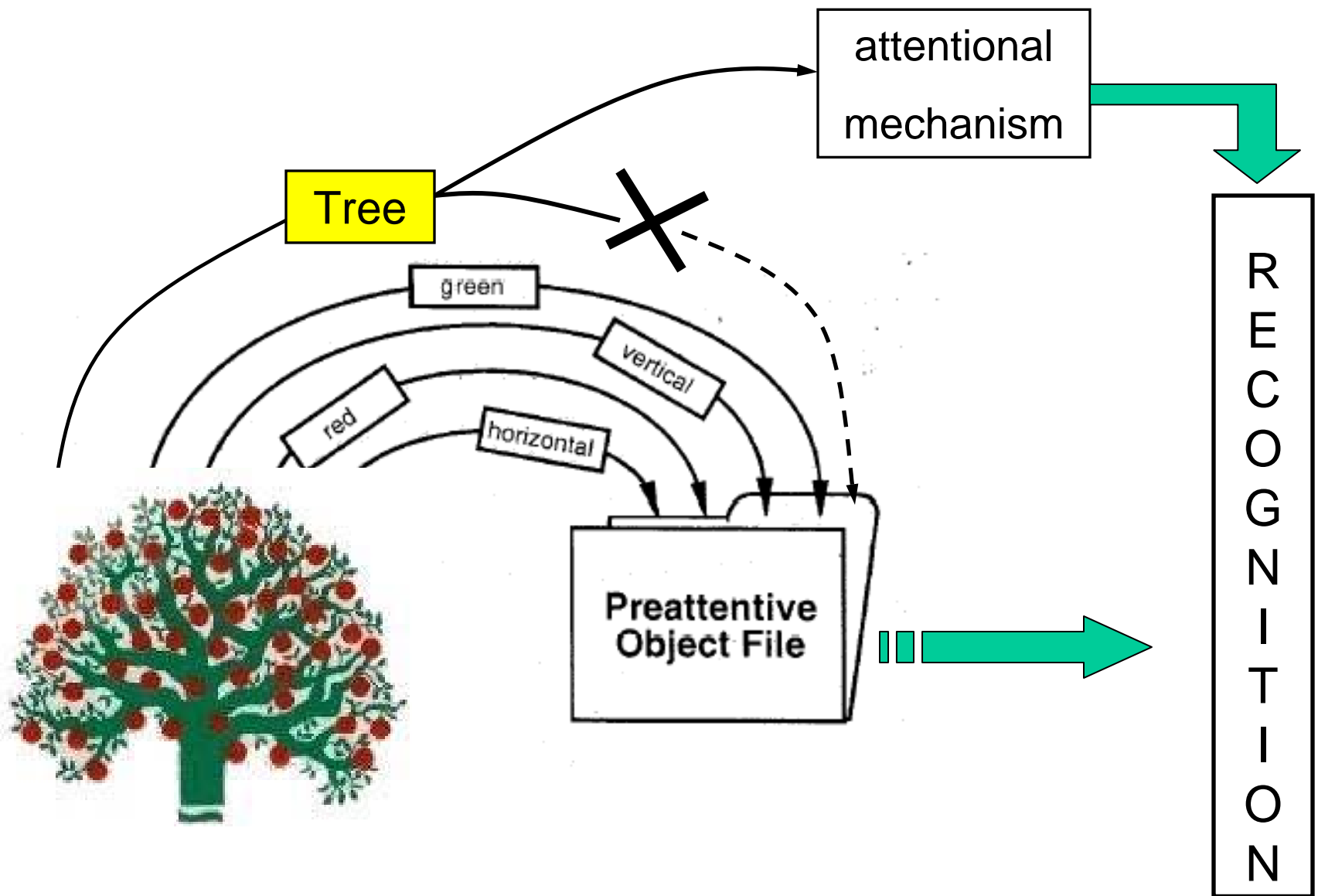
Visual Search:

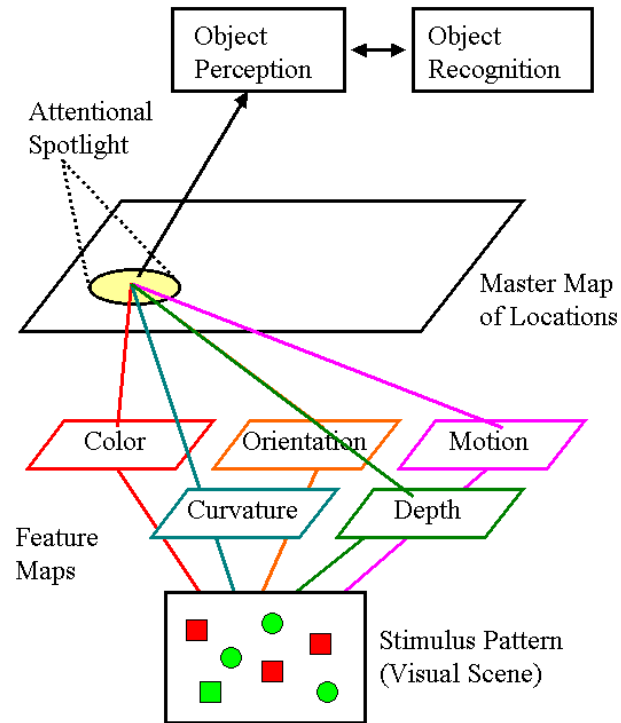
find the green-vertical bar



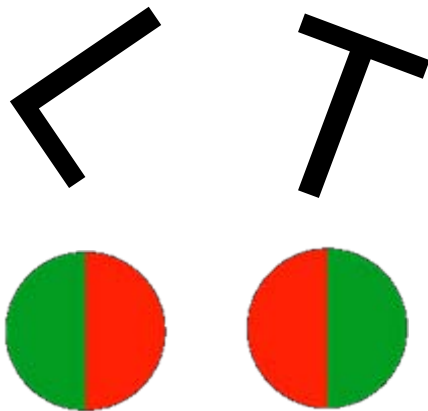
Reaction Time

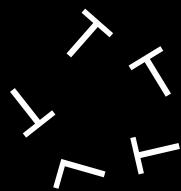


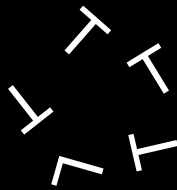


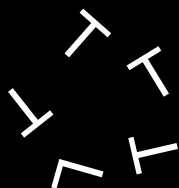


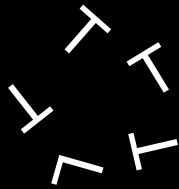
less **attentional load** more



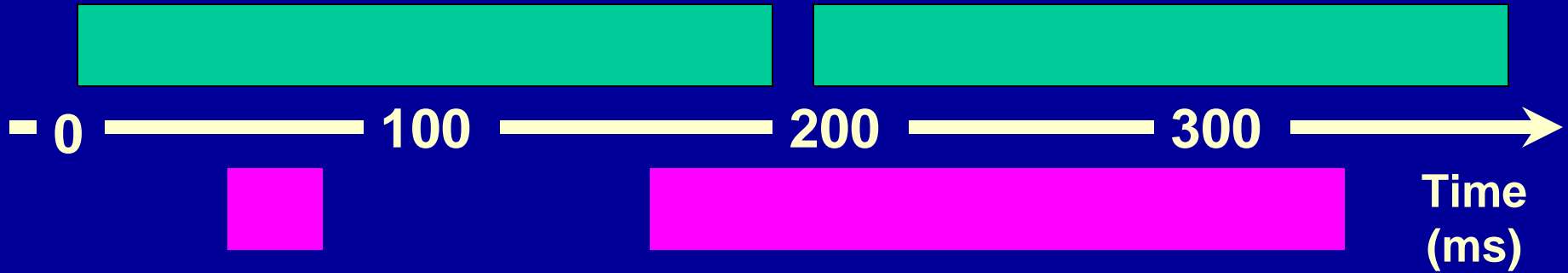
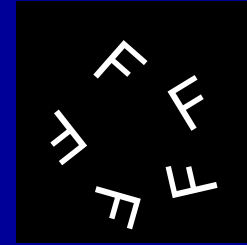
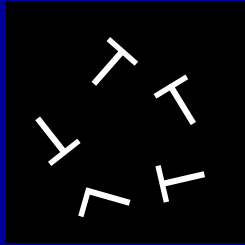




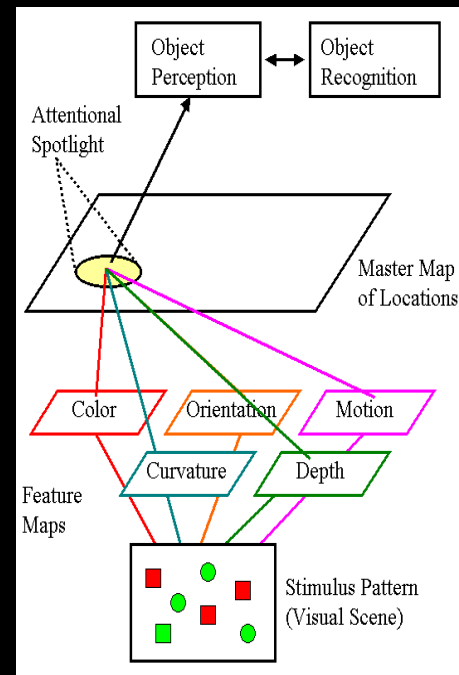
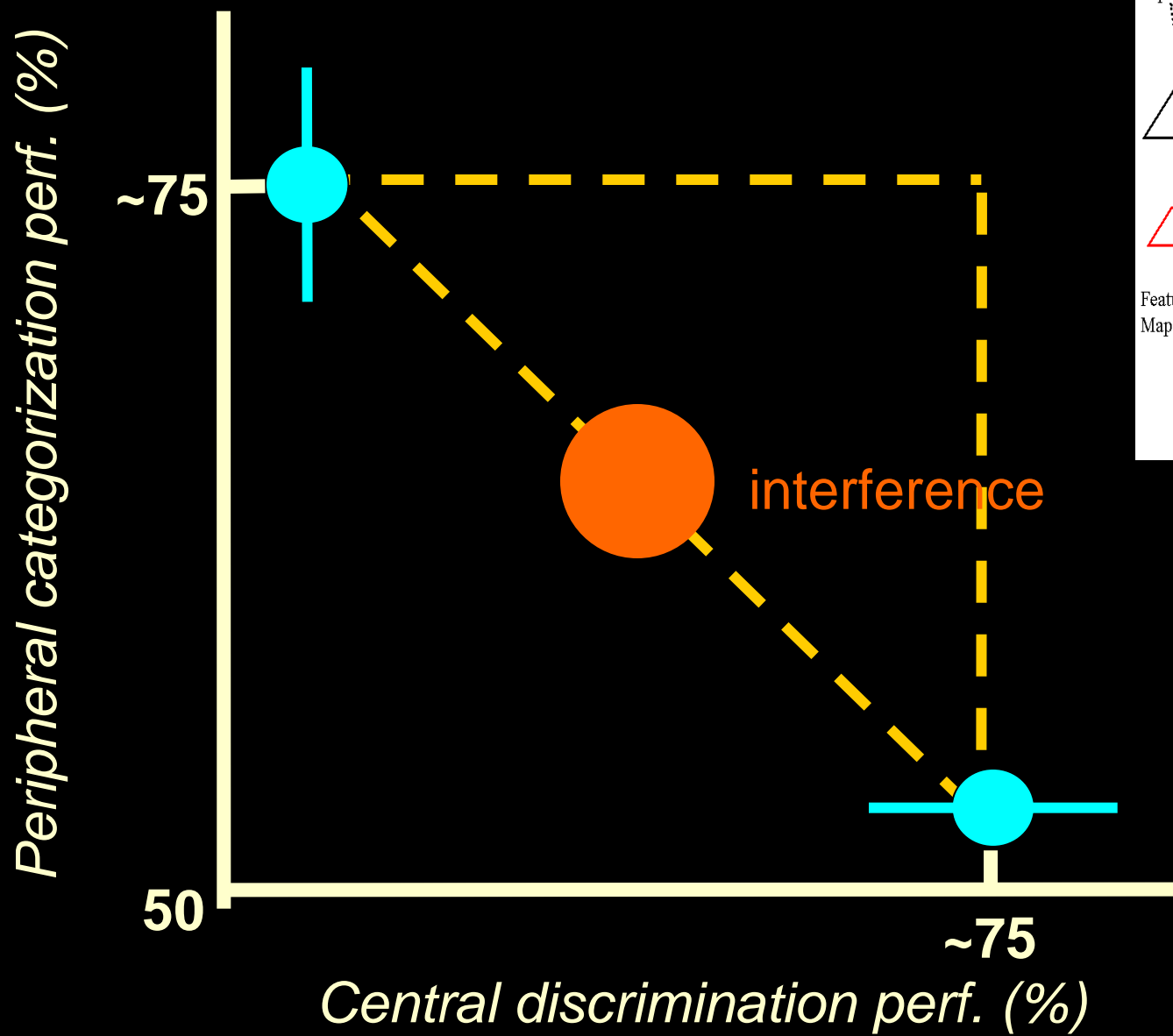




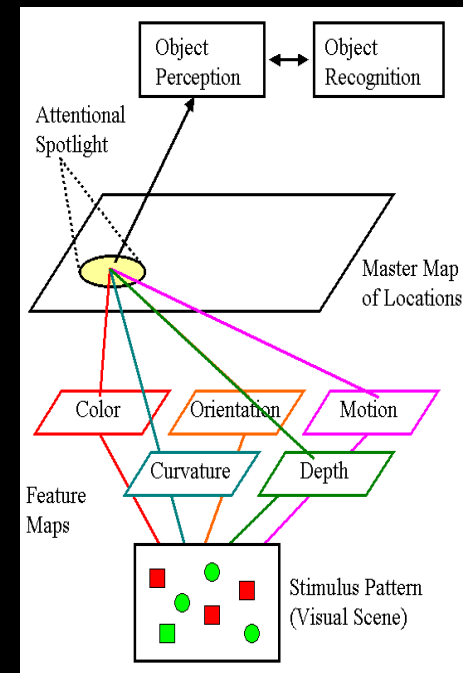
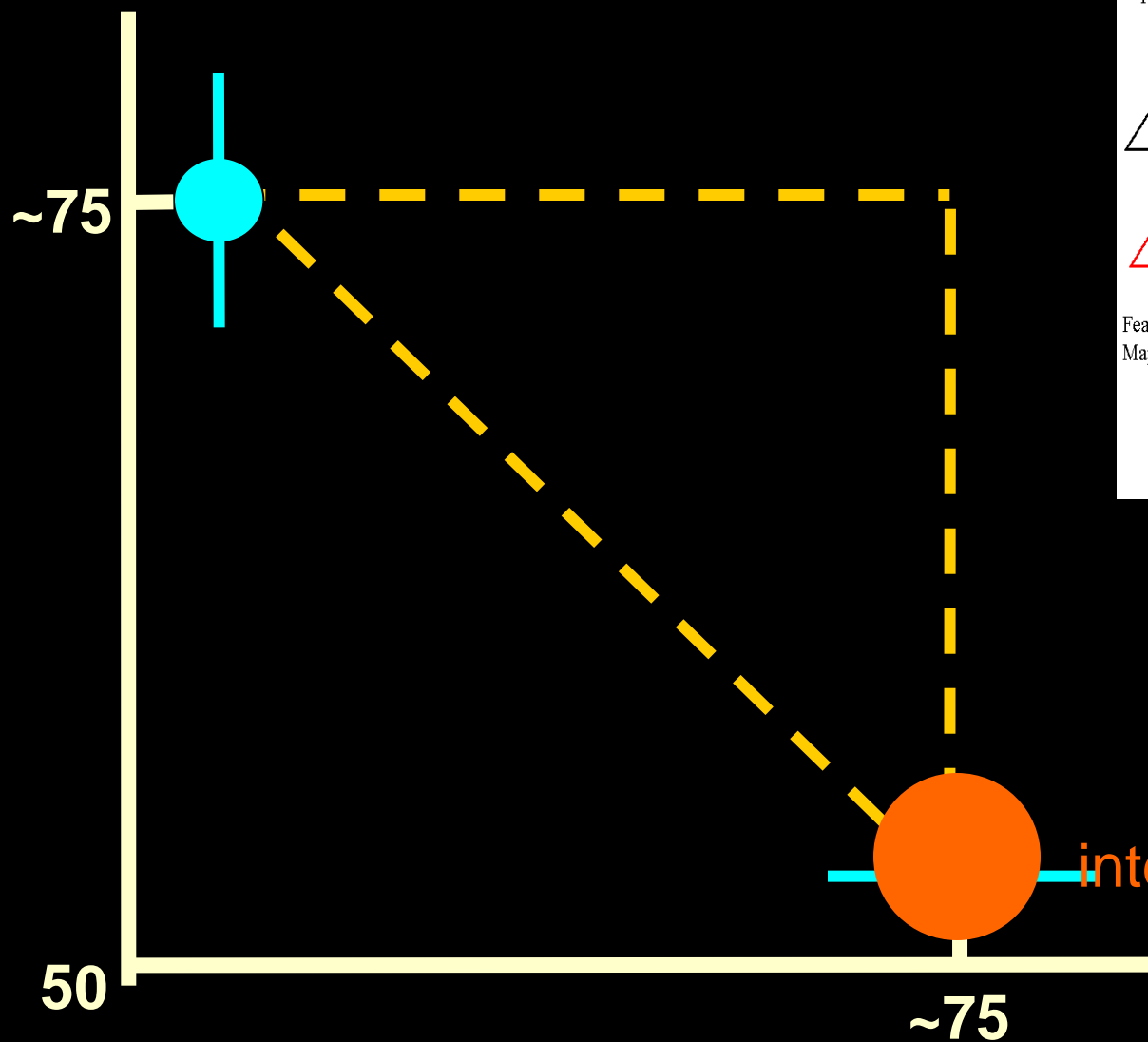
Central

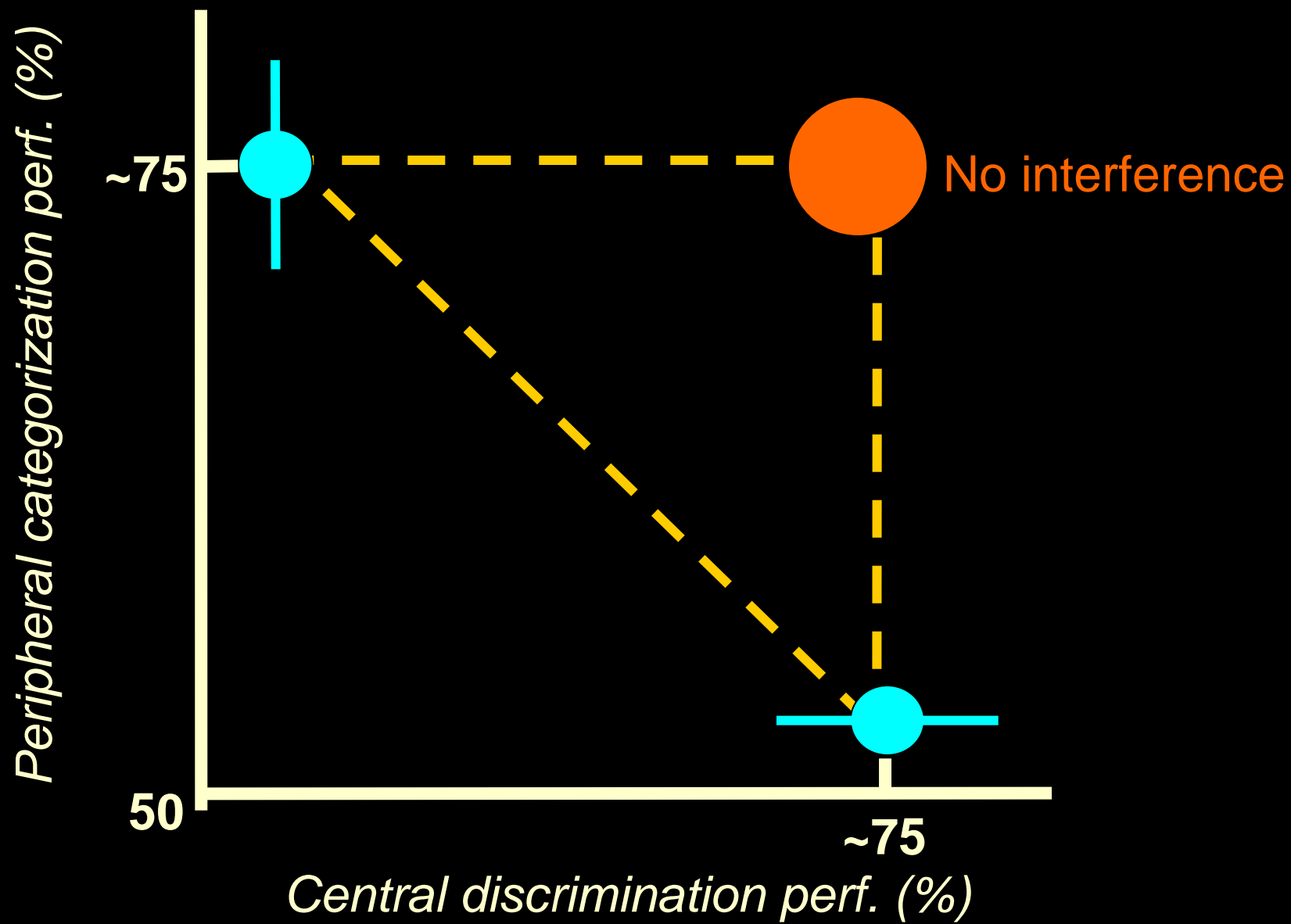


Peripheral



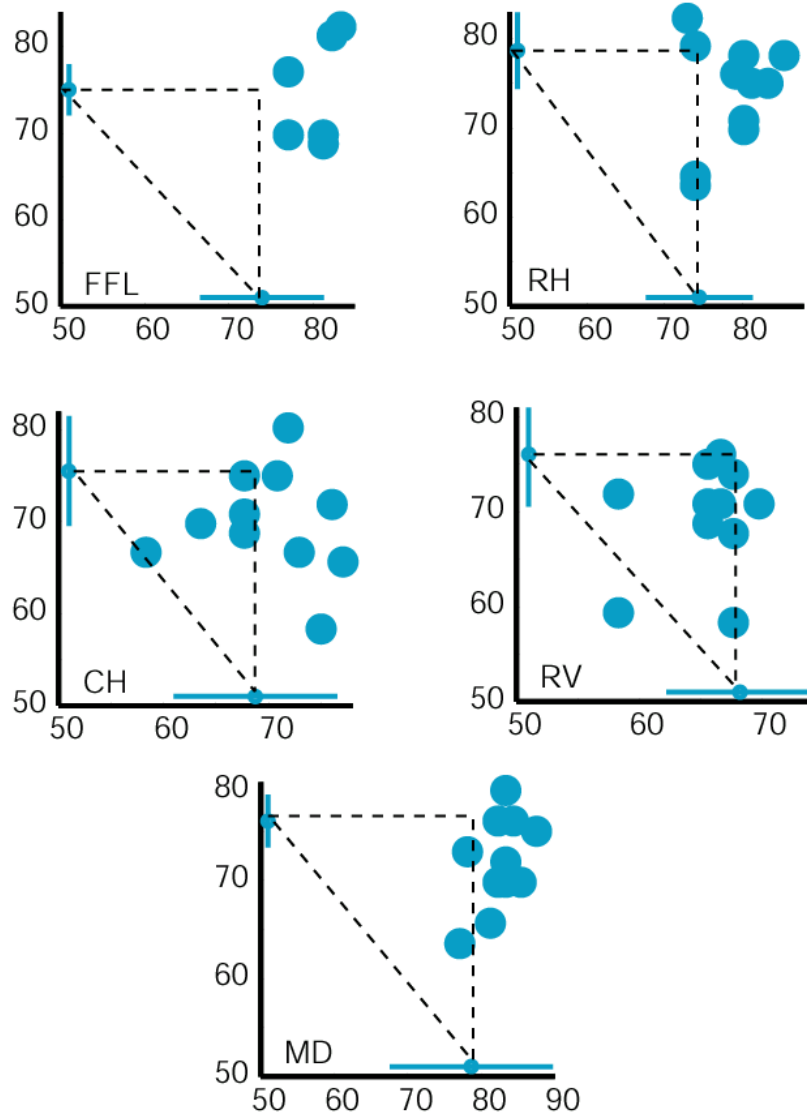
Peripheral categorization perf. (%)



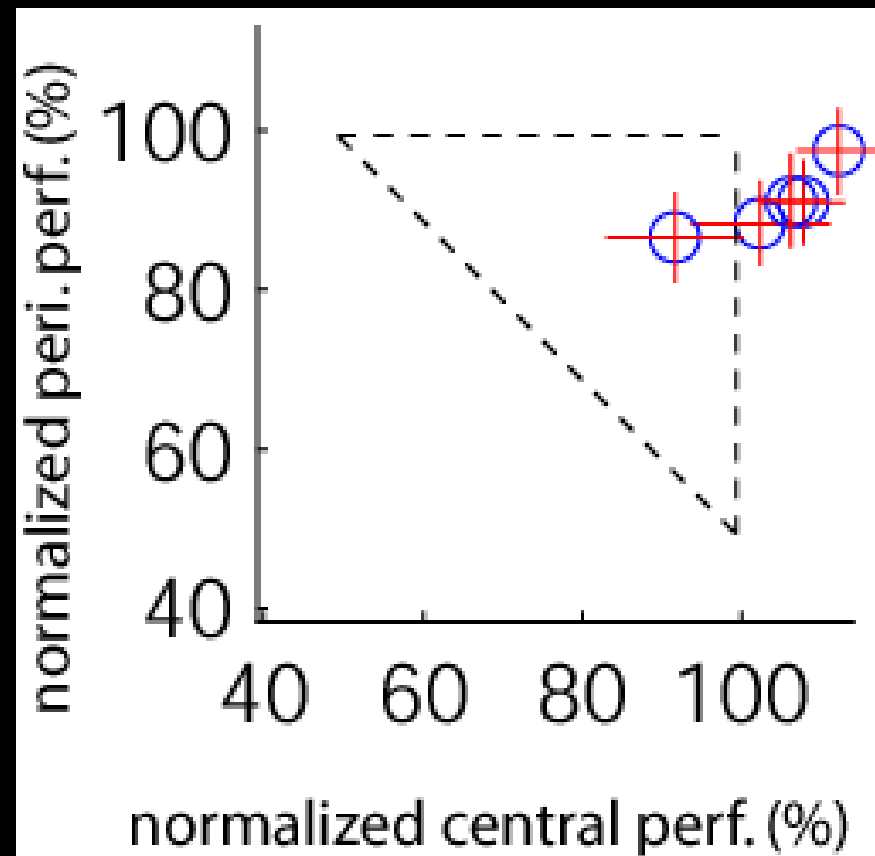


individual results

peripheral task performance (%)

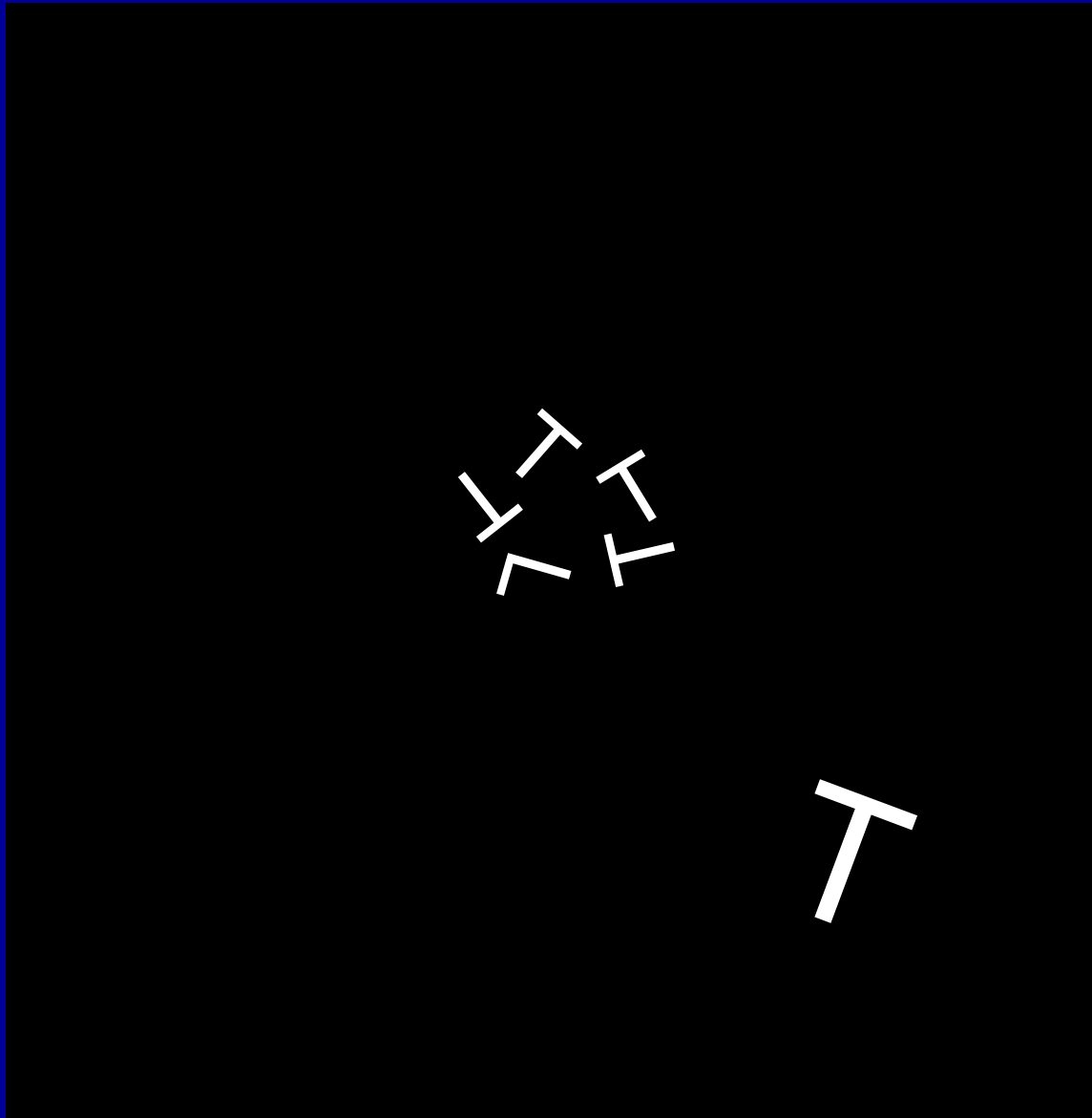


central task performance (%)

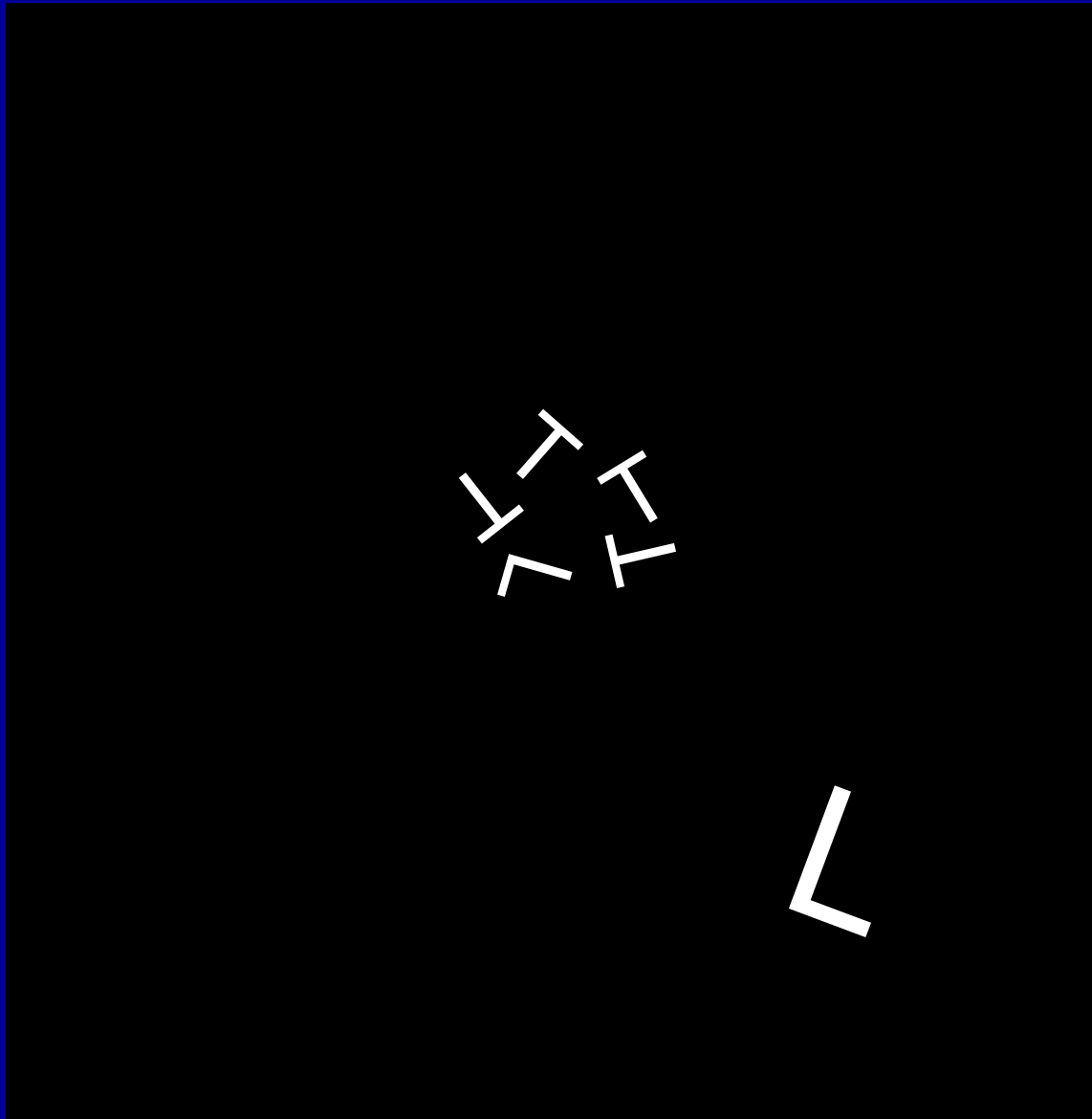


Fei-Fei et al. *PNAS*, 2002

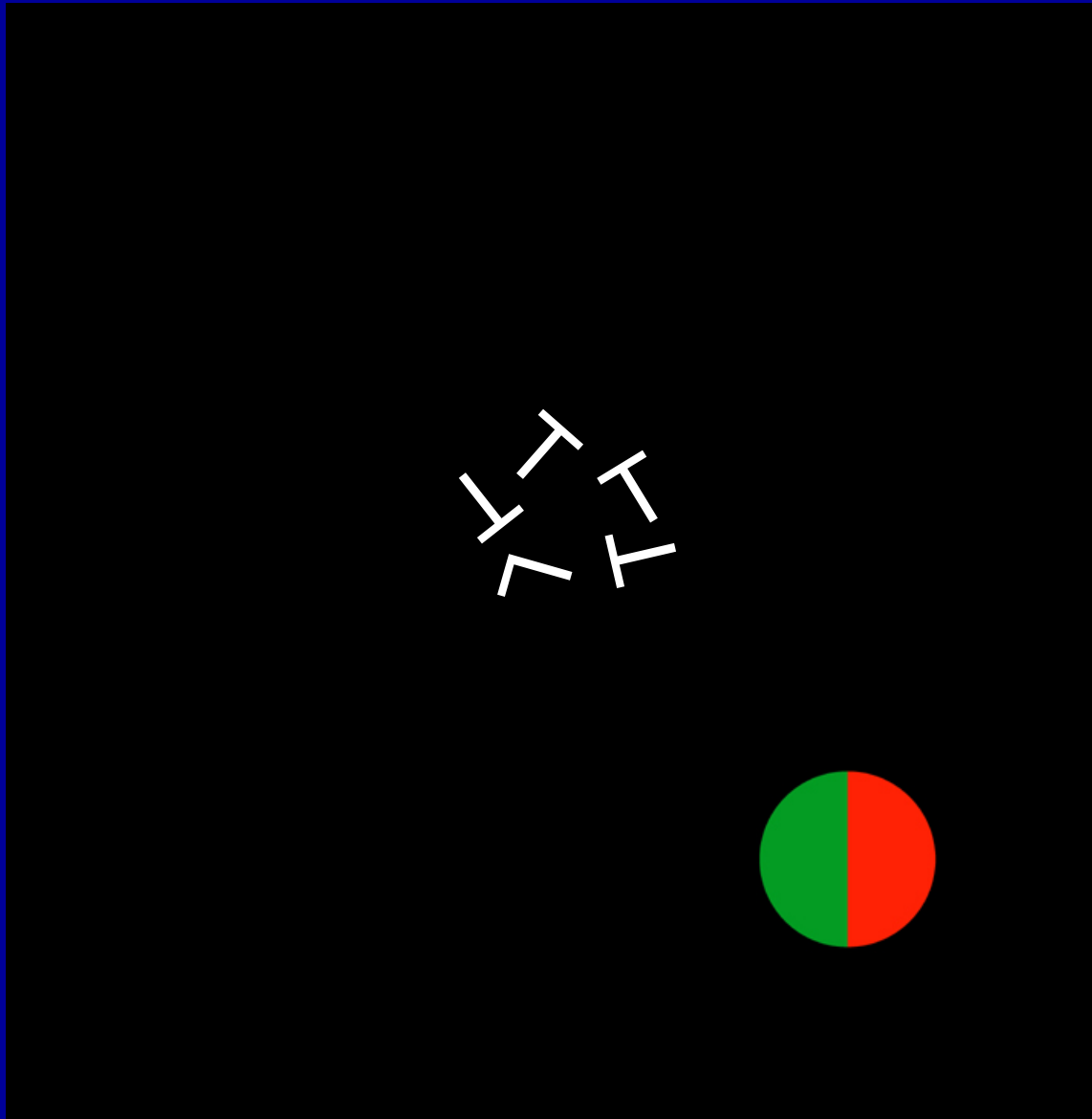
Compare to seemingly simpler tasks



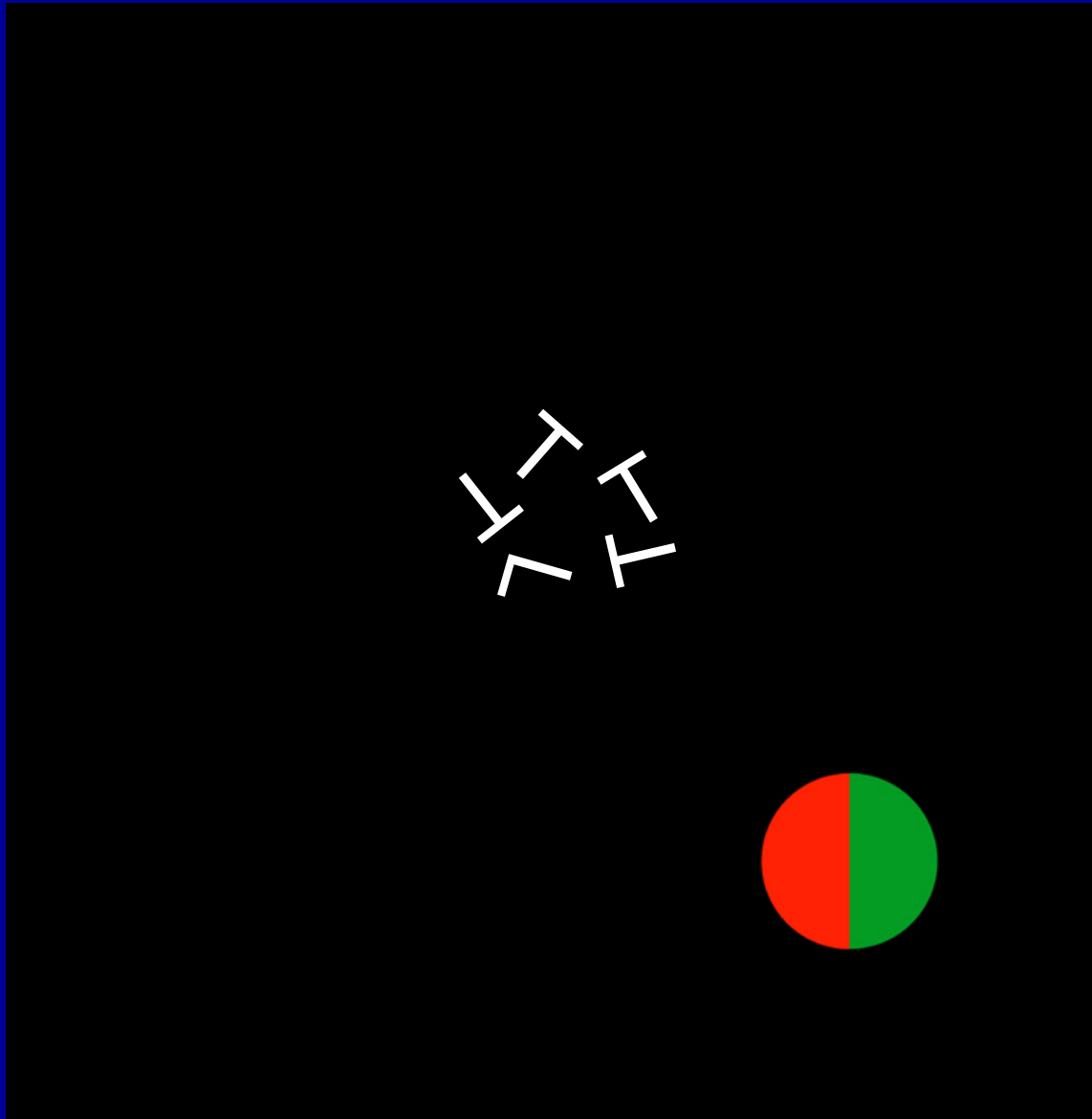
Compare to seemingly simpler tasks



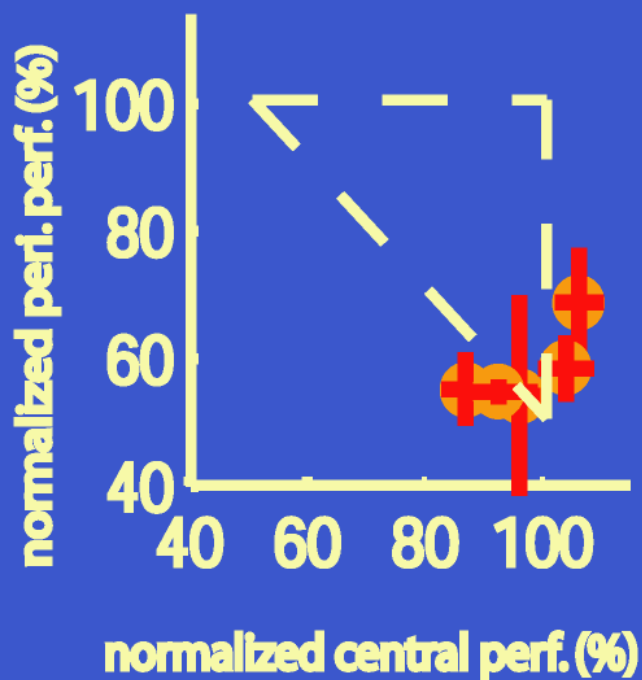
Compare to seemingly simpler tasks



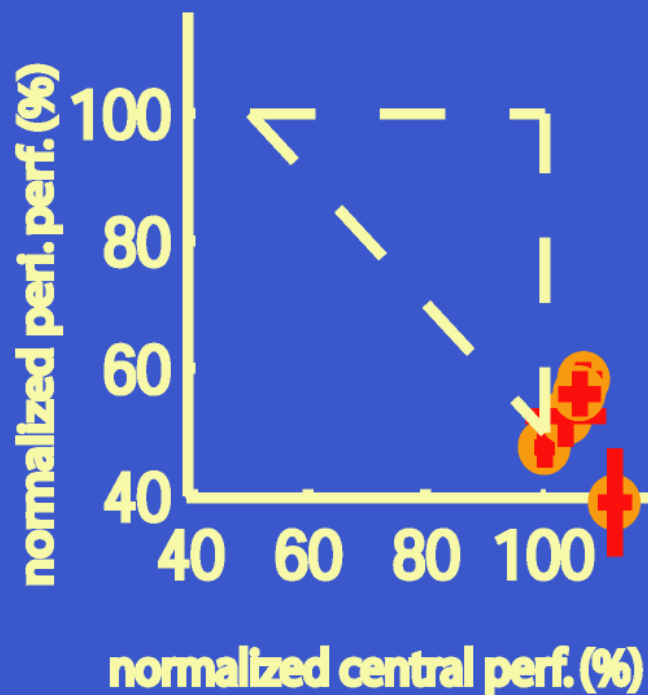
Compare to seemingly simpler tasks



↖ vs. <
(masked by ↗)

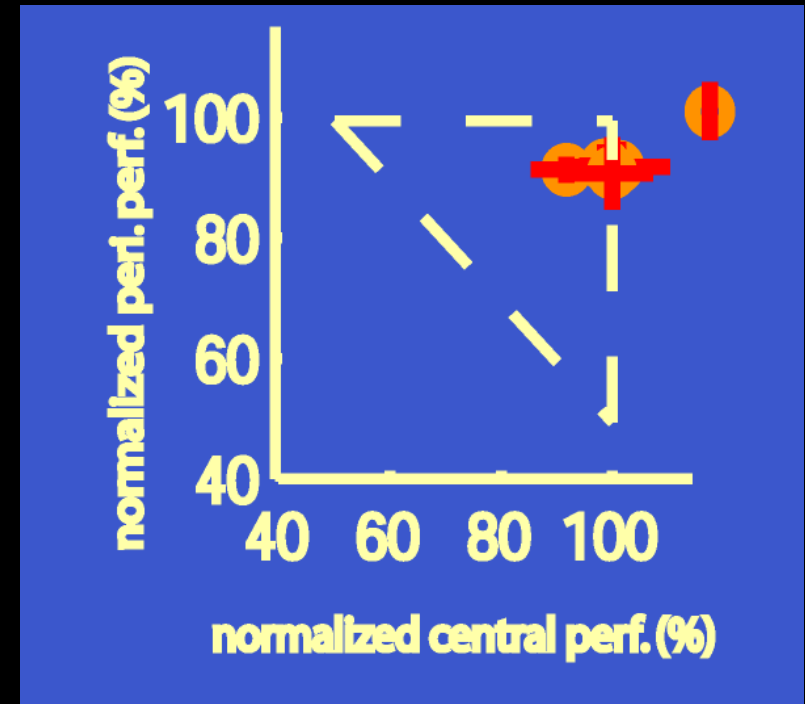
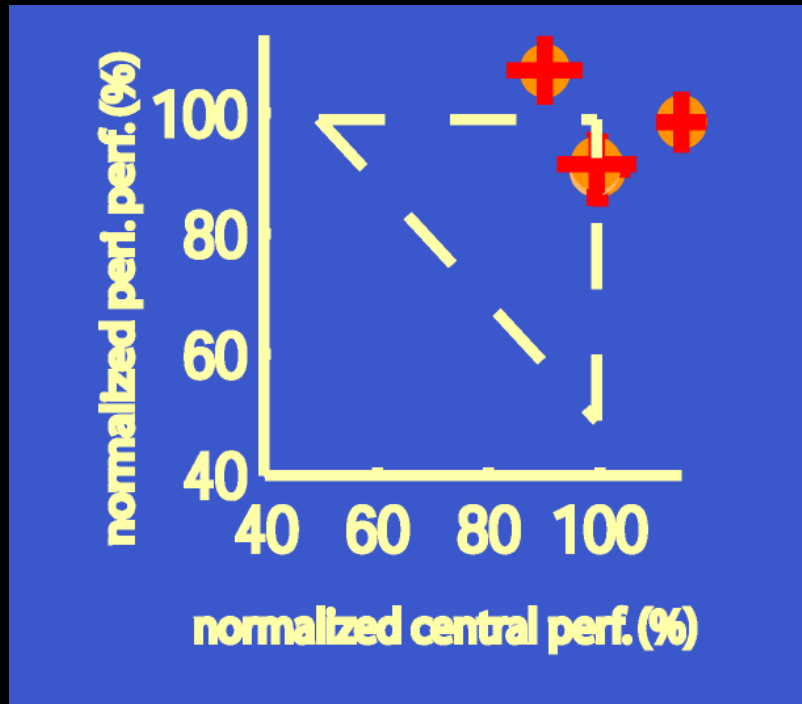


◐ vs. ◑
(masked by ◒)

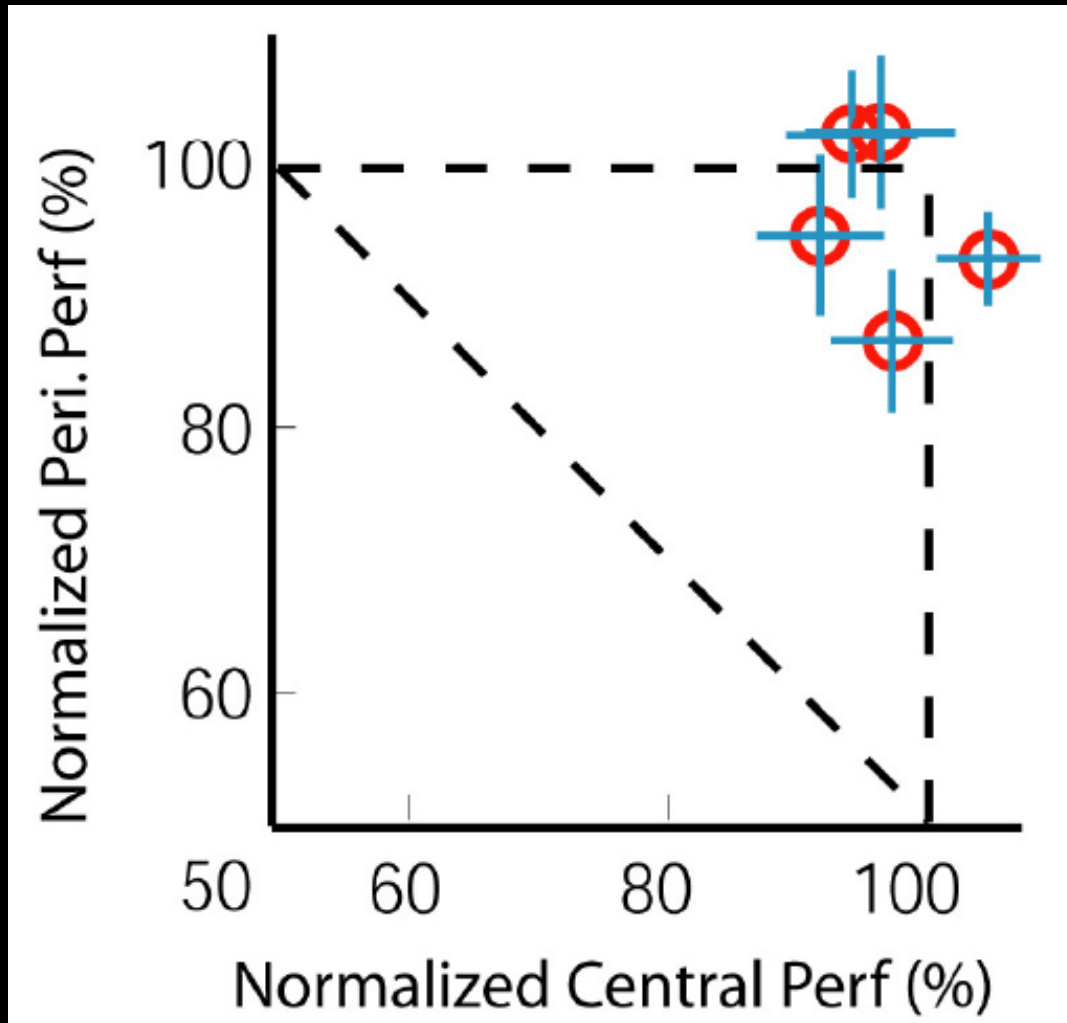


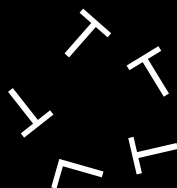
Are animals special?

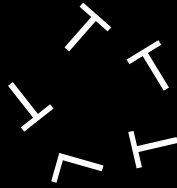


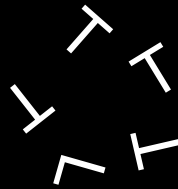


Without color...

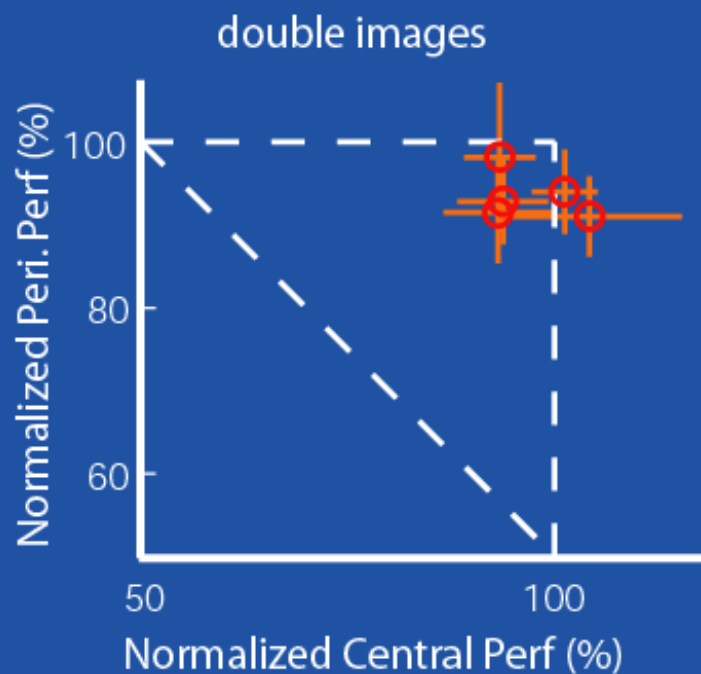
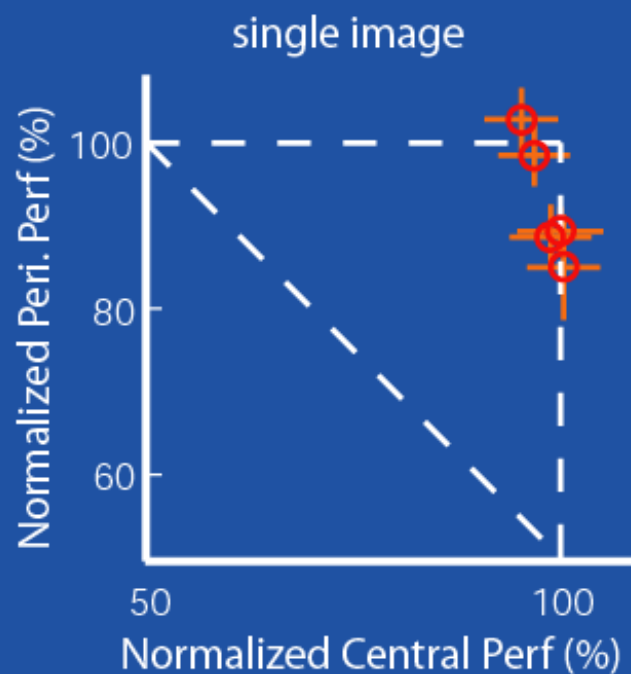


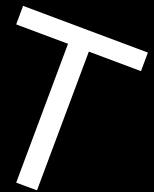
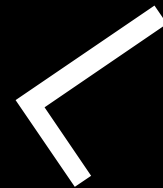






Categorization without attention: Single Image vs. Double Images





Effect of “meaningful” category

randomly rotated


Target

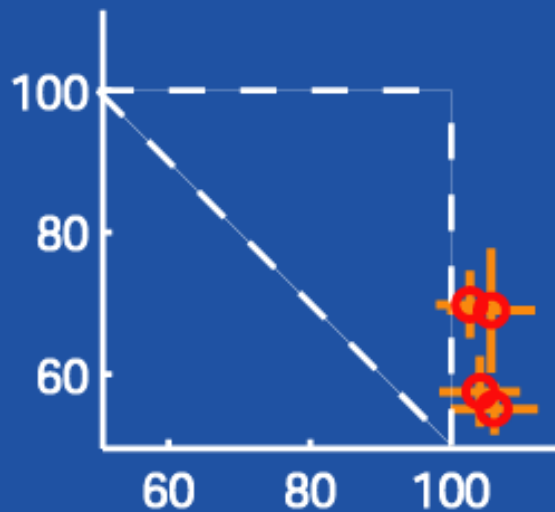
Distractor



vs.



(masked by )



fixed rotation

Target

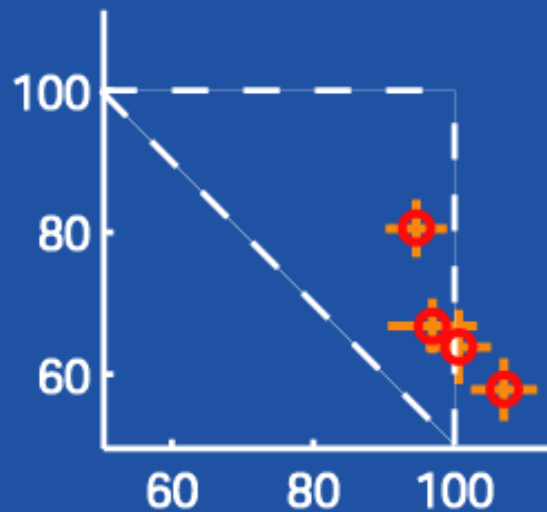
Distractor



vs.



(masked by )



upright position


Target

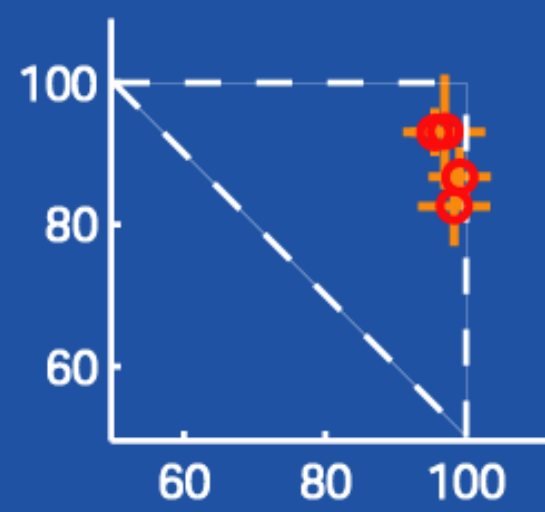
Distractor



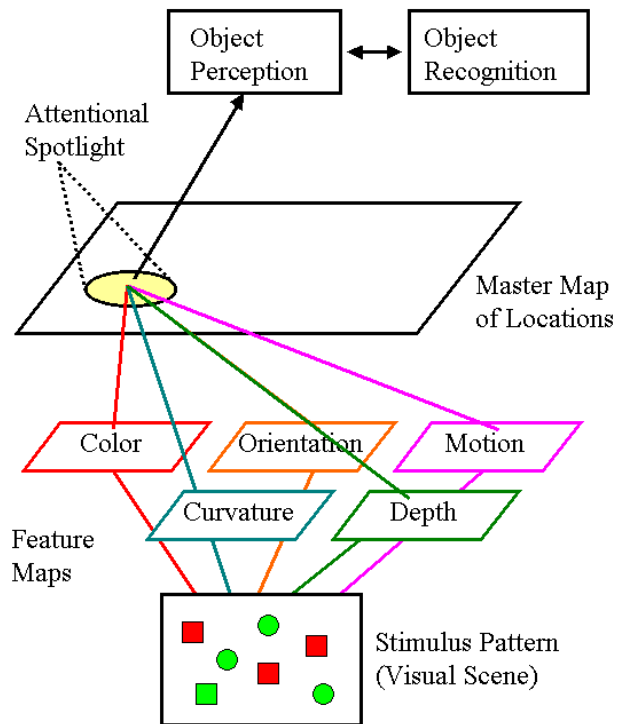
vs.



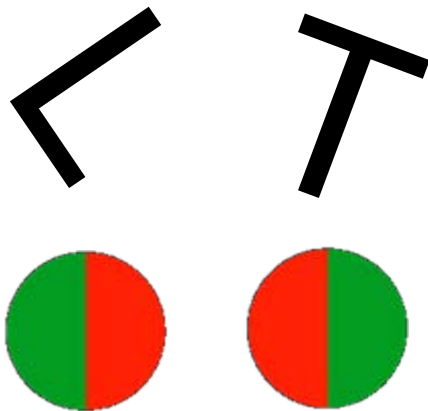
(masked by )



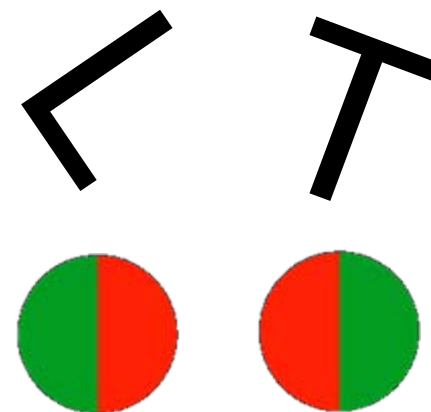
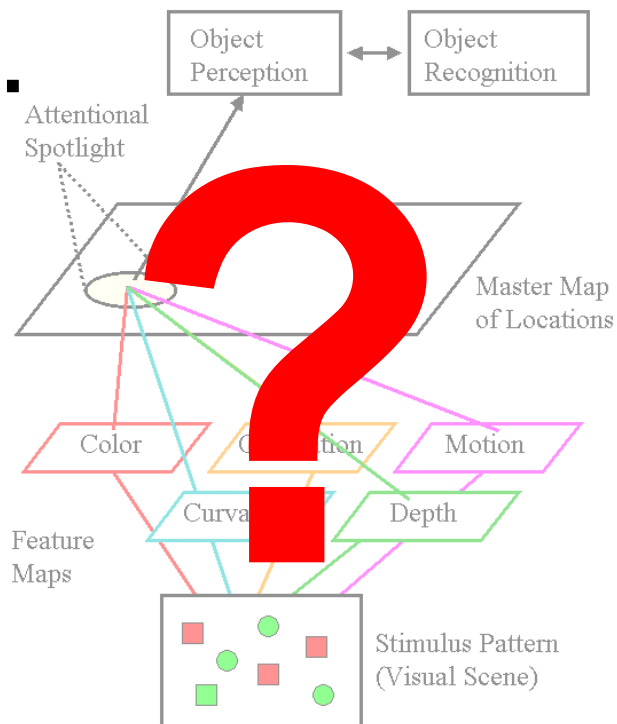
F.I.T. predicted...



less **attentional load** more

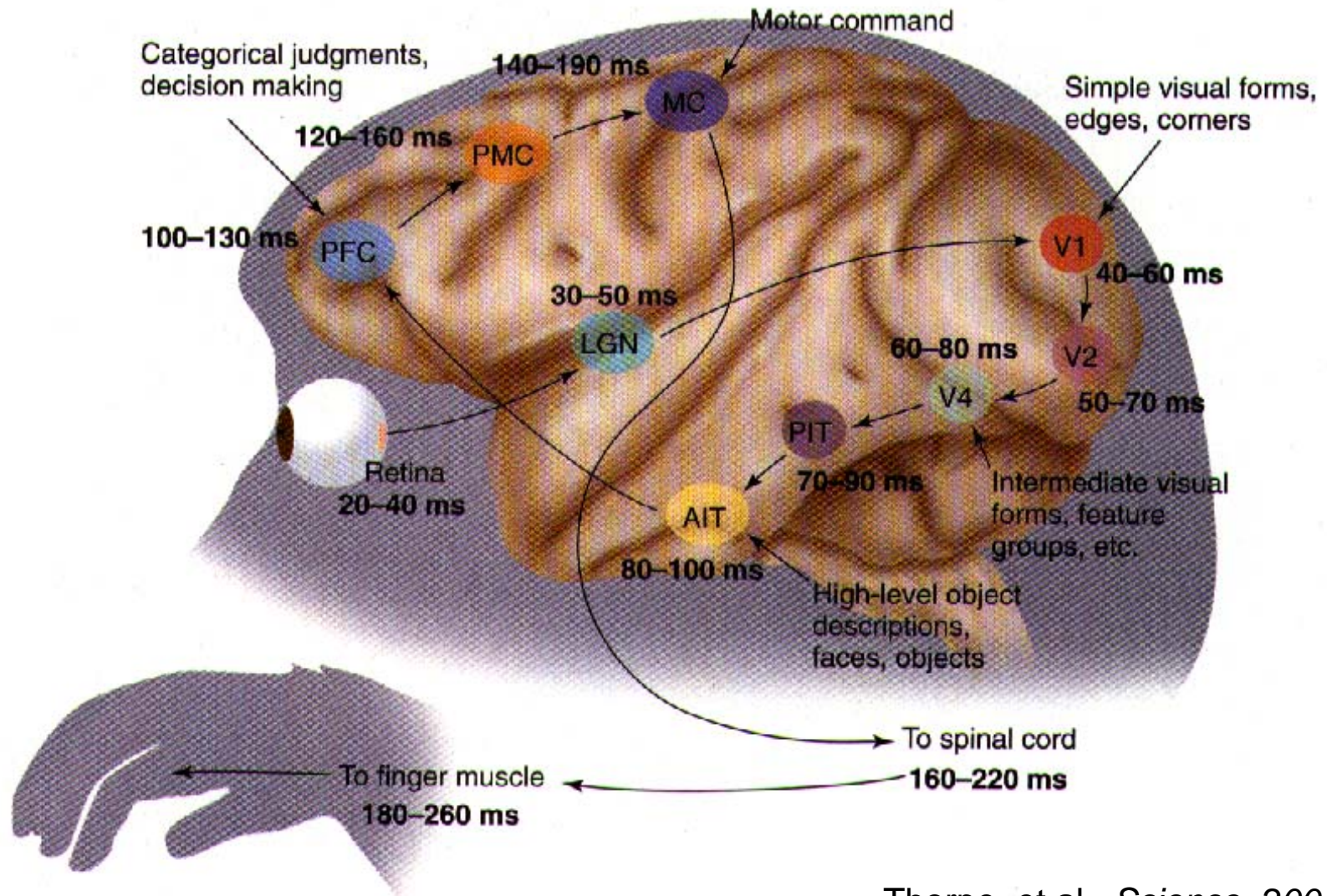


Our data shows...



Rapid Perception of Natural Scenes

- Where/how does this happen?



Beaches



Highways



Forests



Buildings

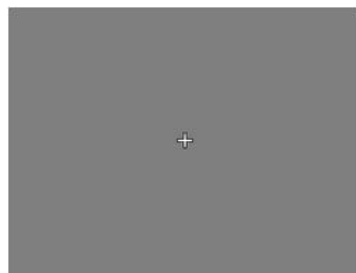


Mountains



Industry

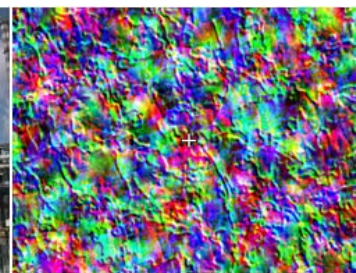




500 ms



32-45 ms

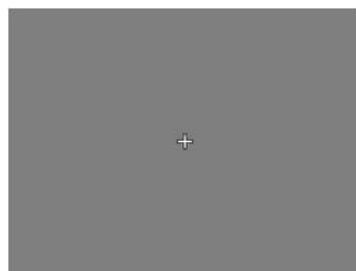


500 ms



< 2000 ms

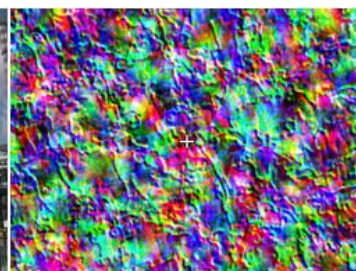




500 ms



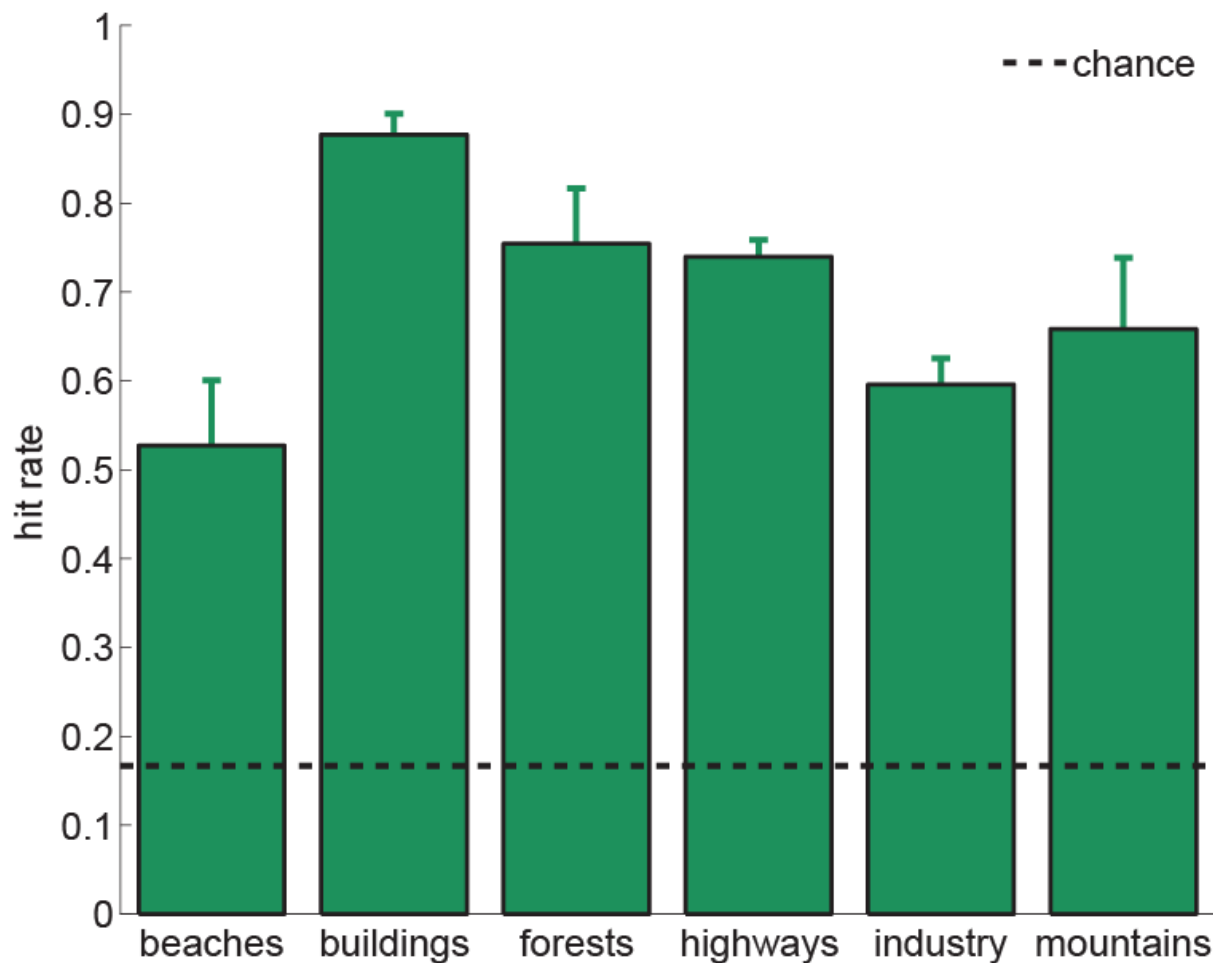
32-45 ms



500 ms



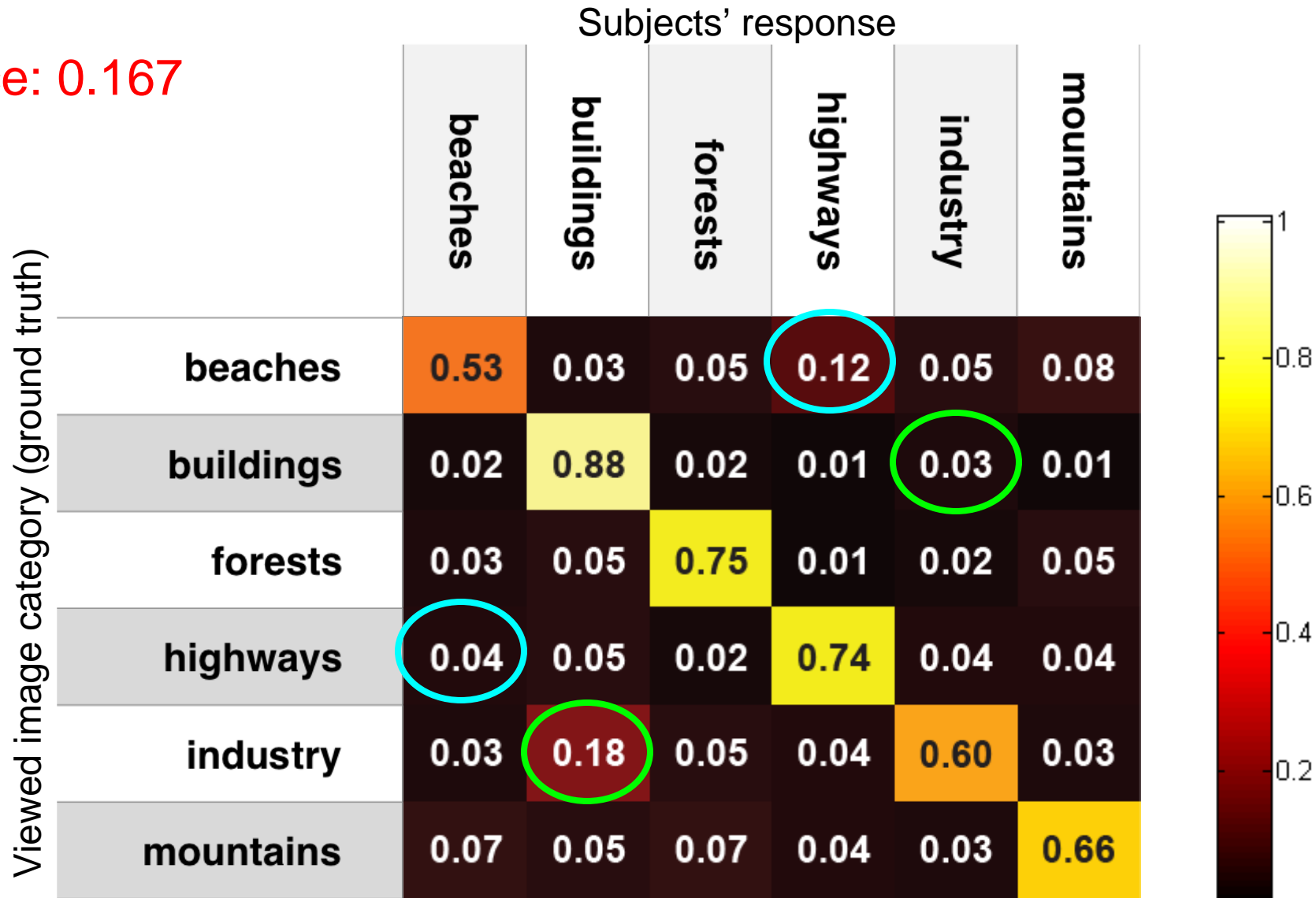
< 2000 ms



6 AFC, N = 4, error bars: s.e.m.

Behavioral Performance

chance: 0.167



Beaches



Highways



Forests



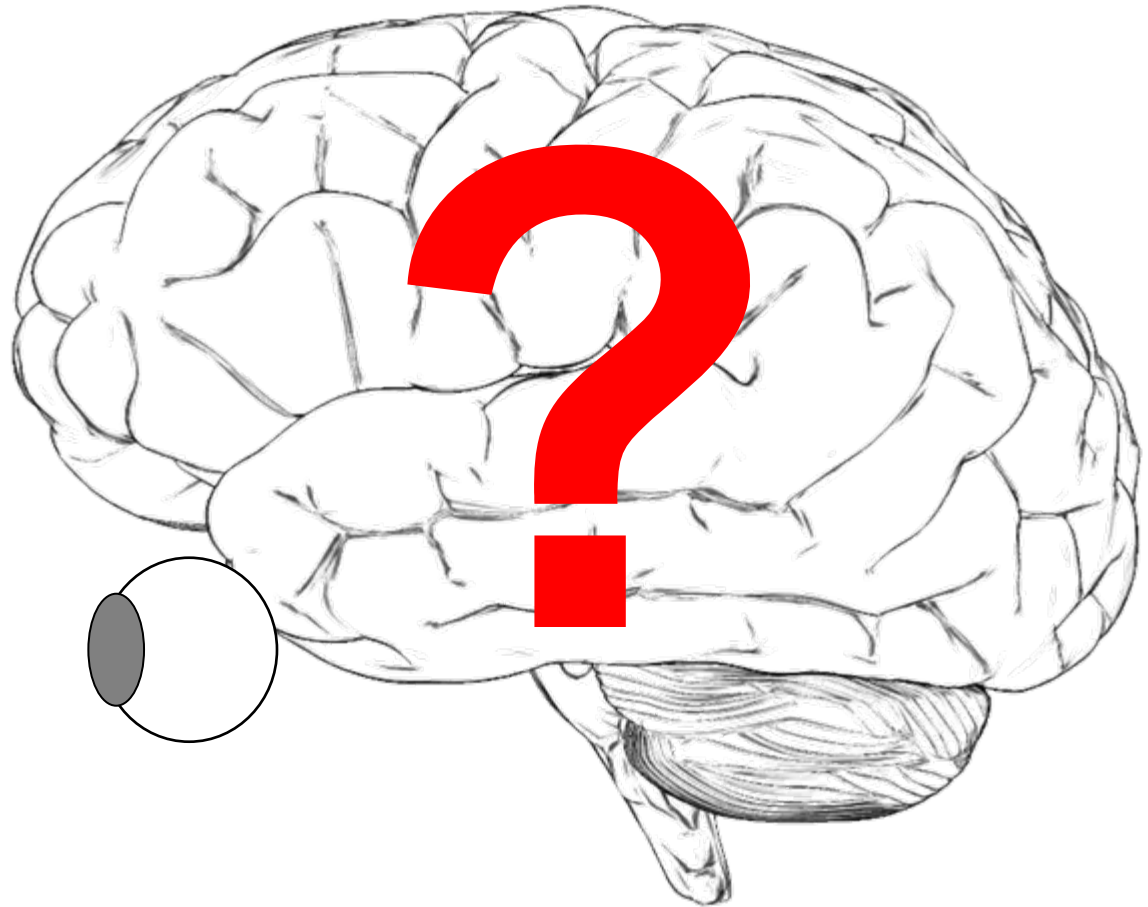
Buildings



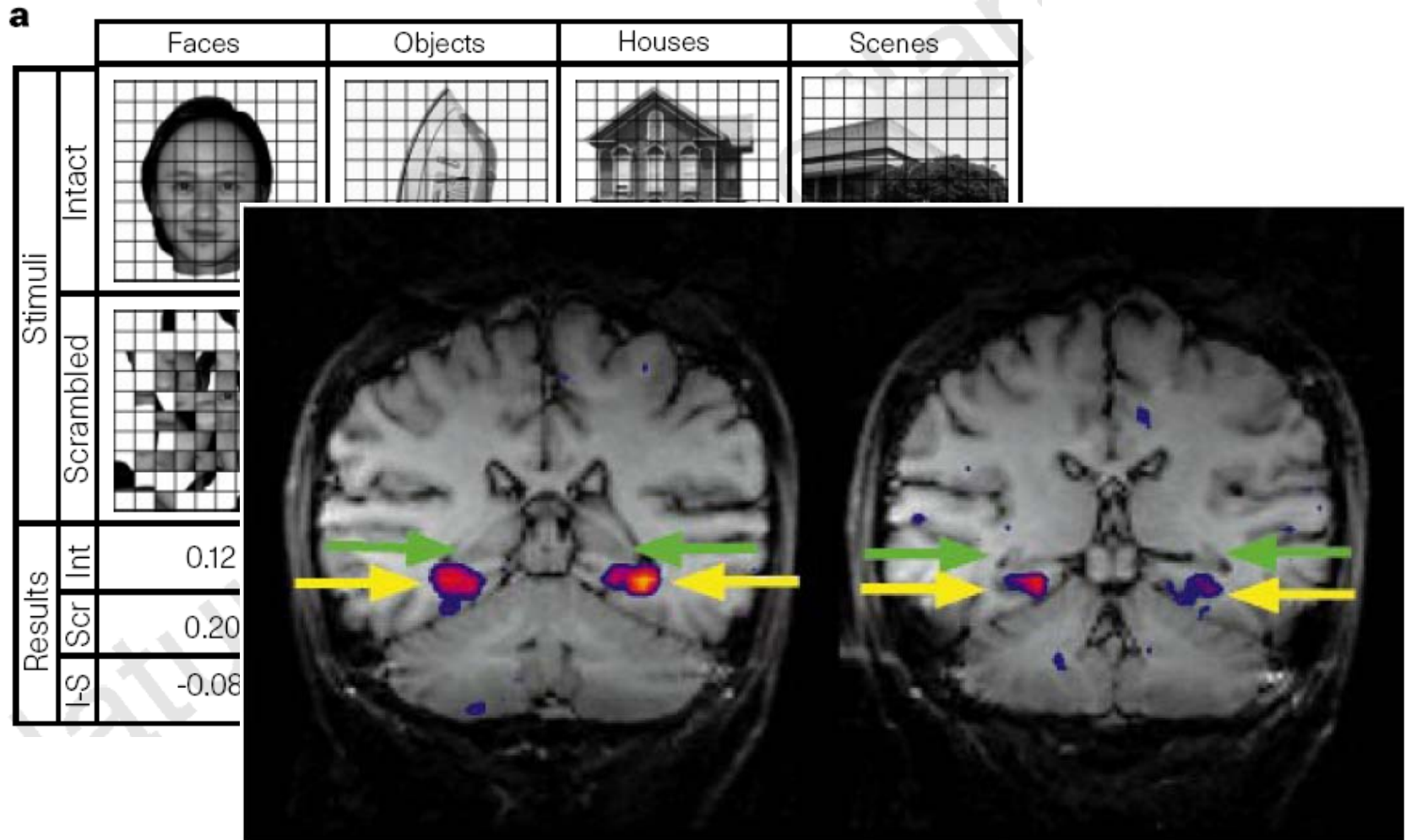
Mountains



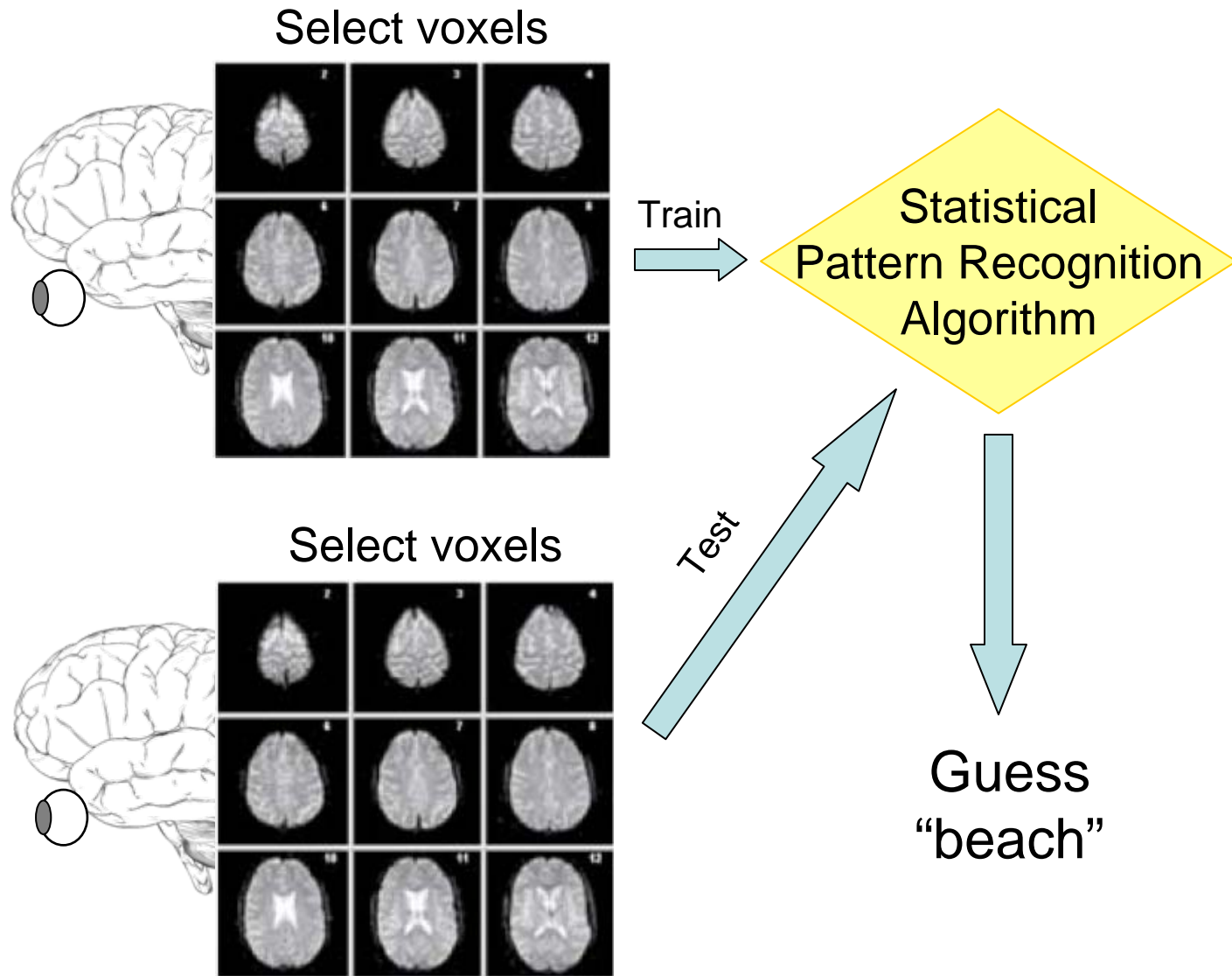
Industry



PPA: Parahippocampal Place Area

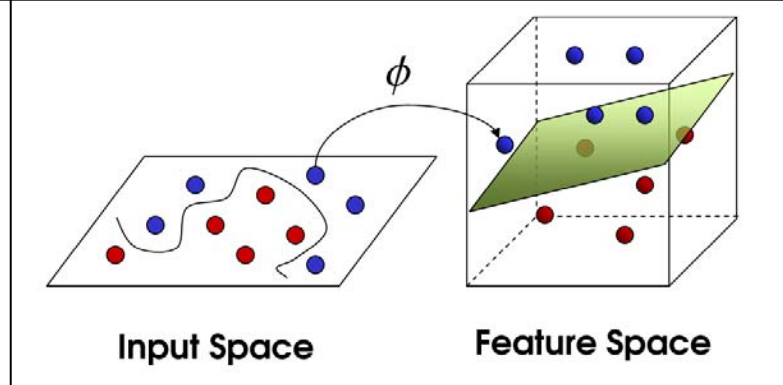


Pattern Recognition

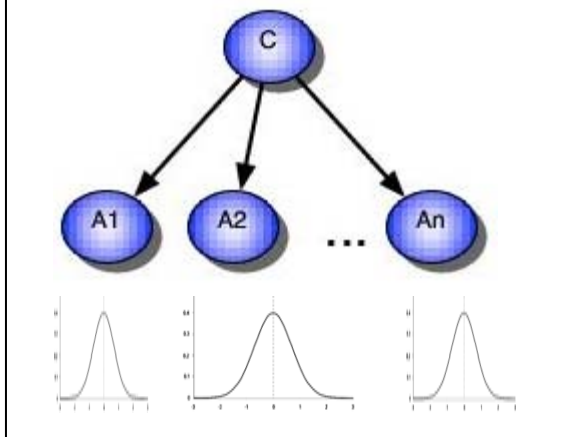


Pattern Recognition

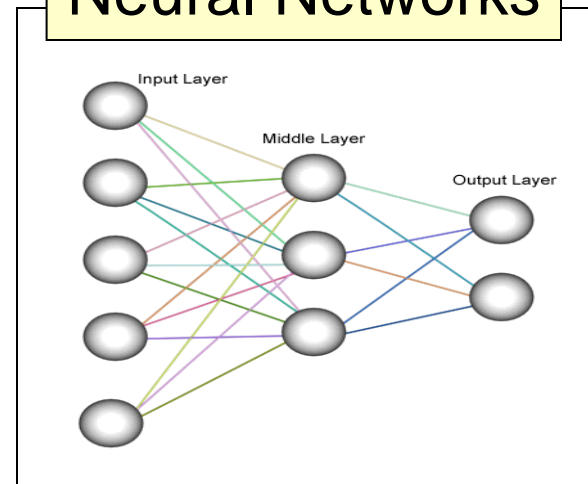
Support Vector Machine (SVM)



Gaussian Naïve Bayes (GNB)

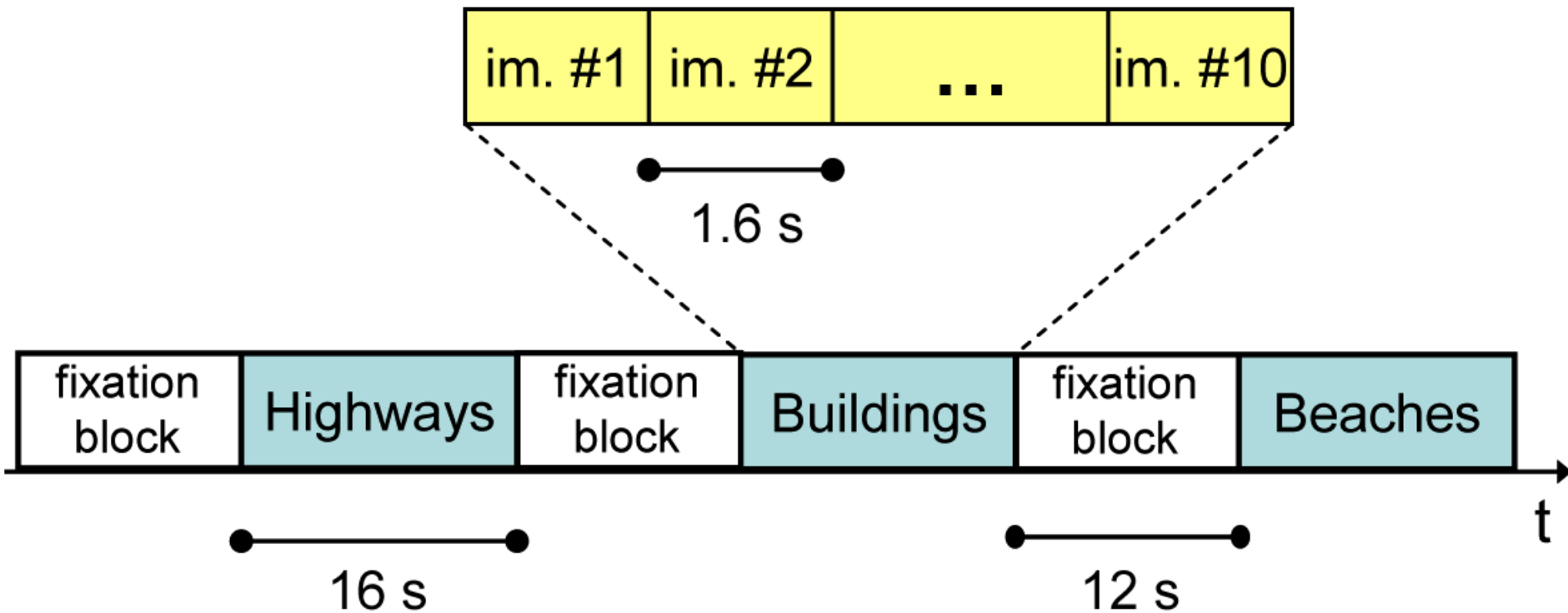


Neural Networks



Statistical
Pattern Recognition
Algorithm

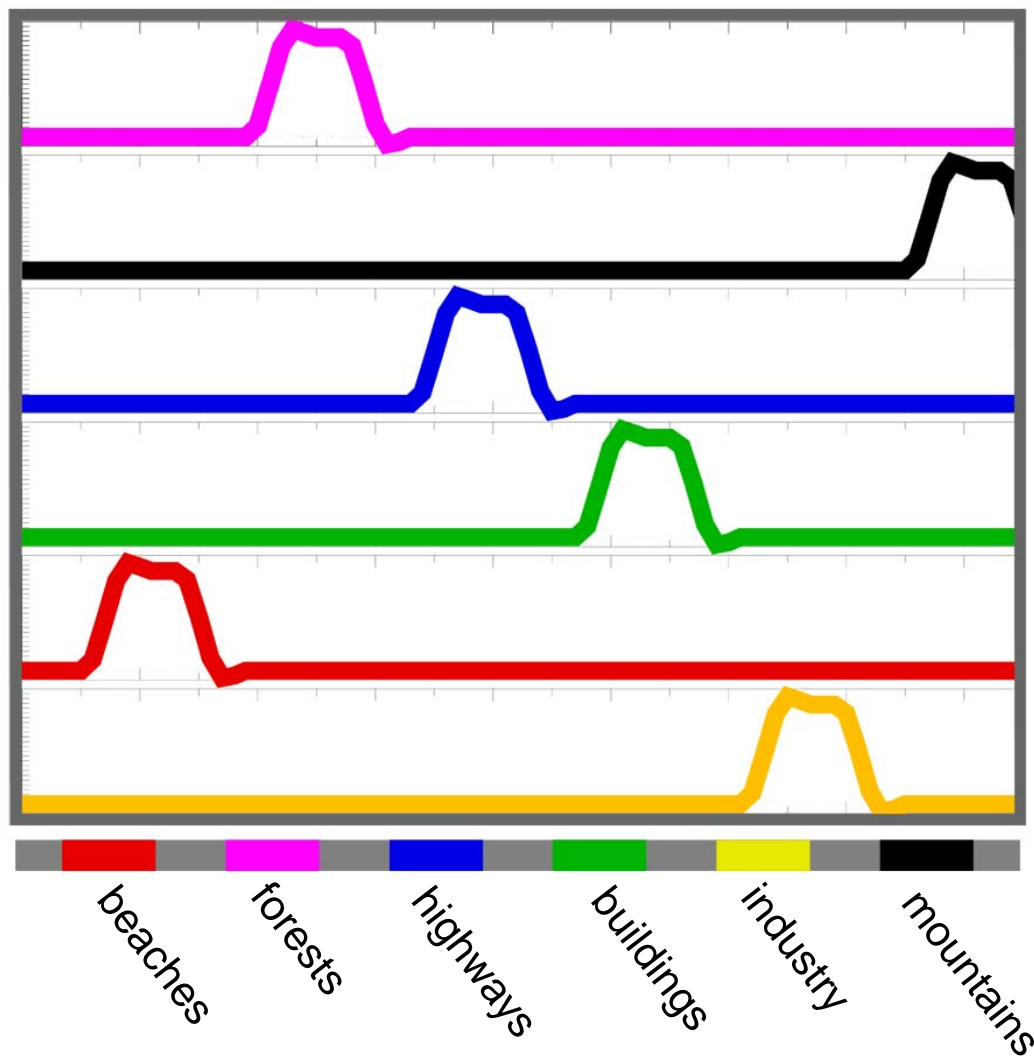
Experimental Setup (fMRI)



- 6 blocks per run (all 6 categories)
- 12 runs for each subject
- Alternating runs feature upright or inverted images

Voxel Selection

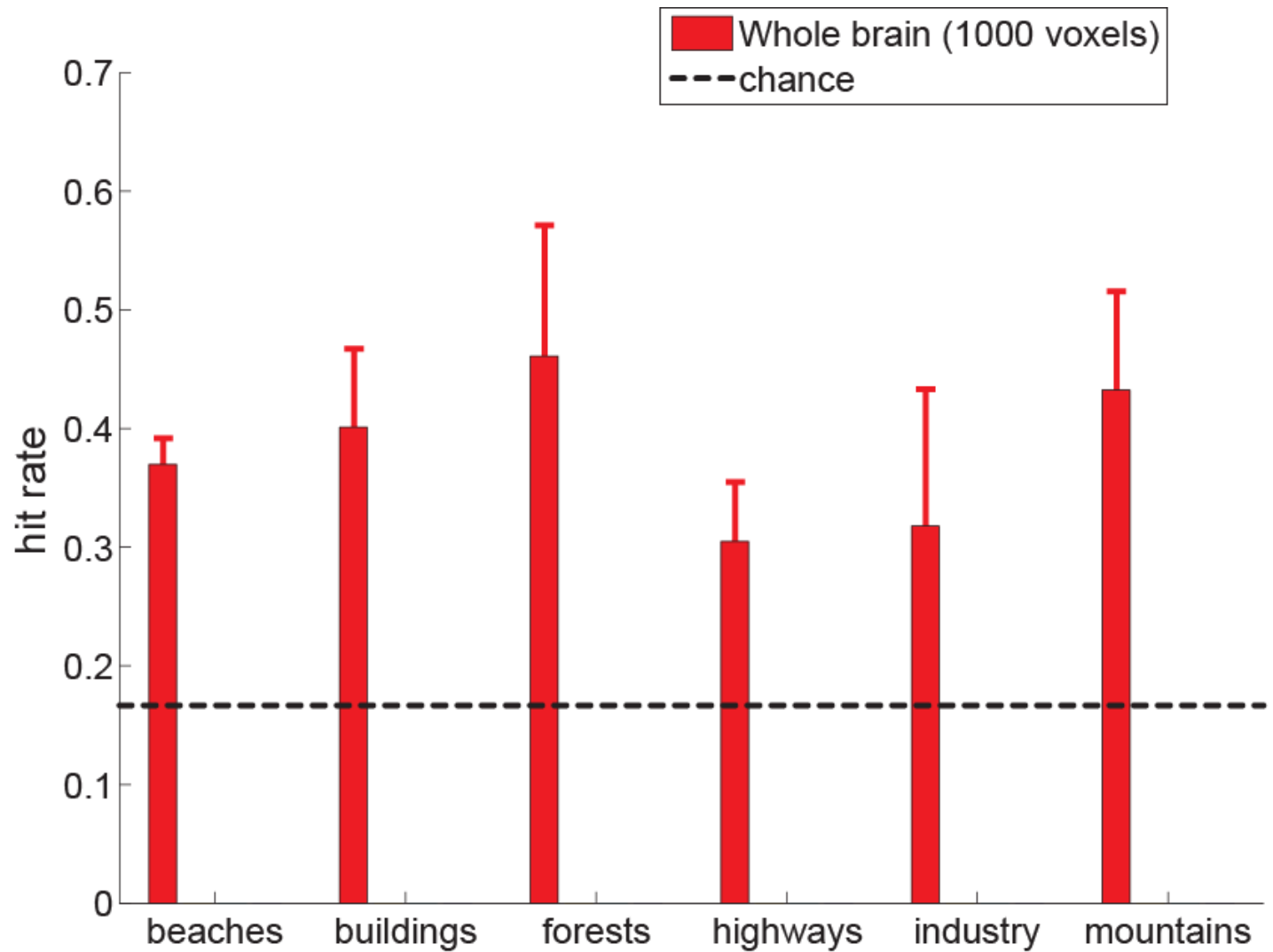
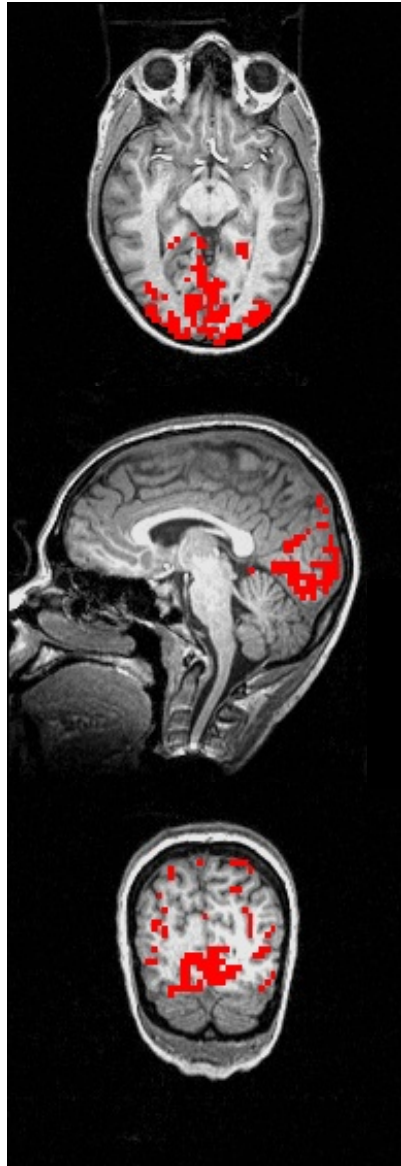
Univariate Multiple Regression



F-statistic

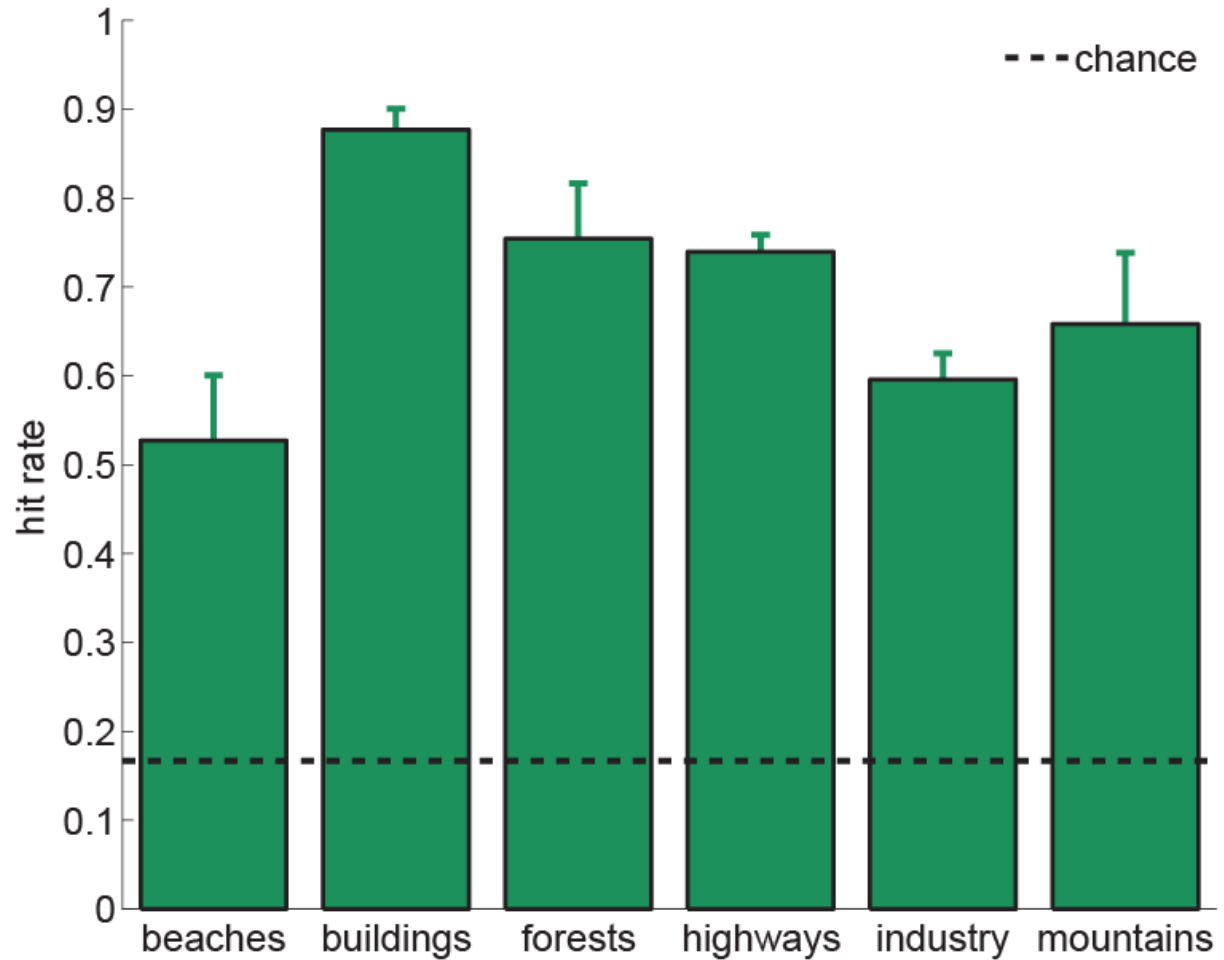
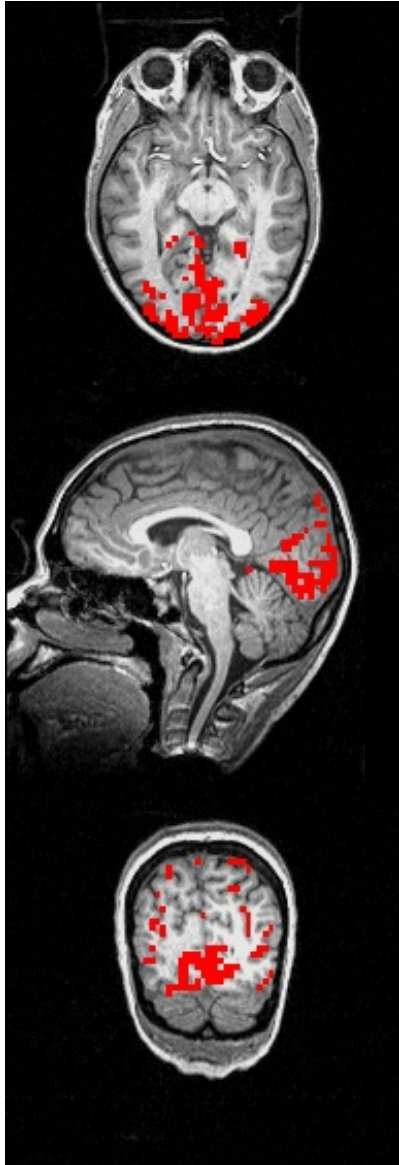


Decoding Performance

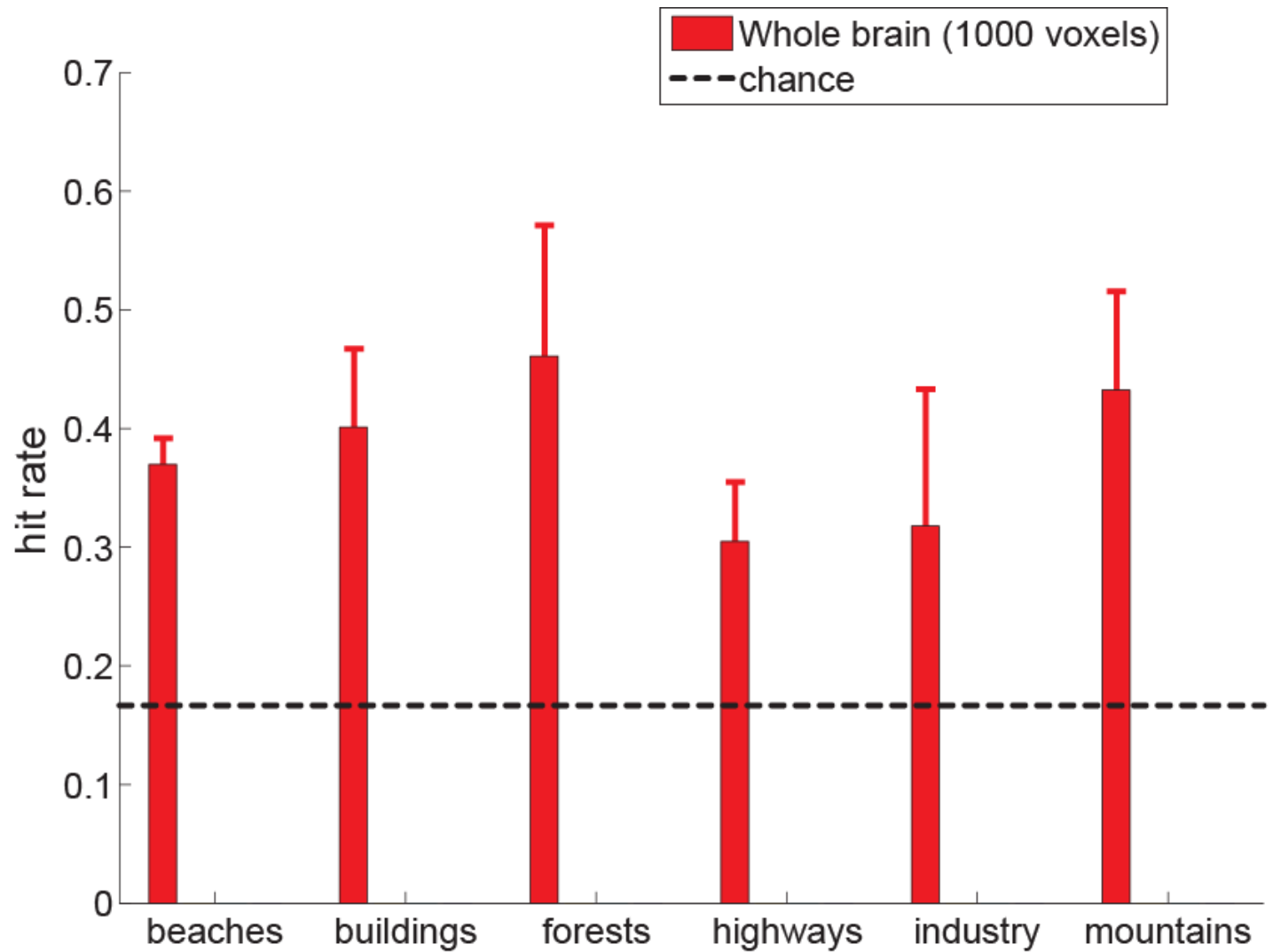
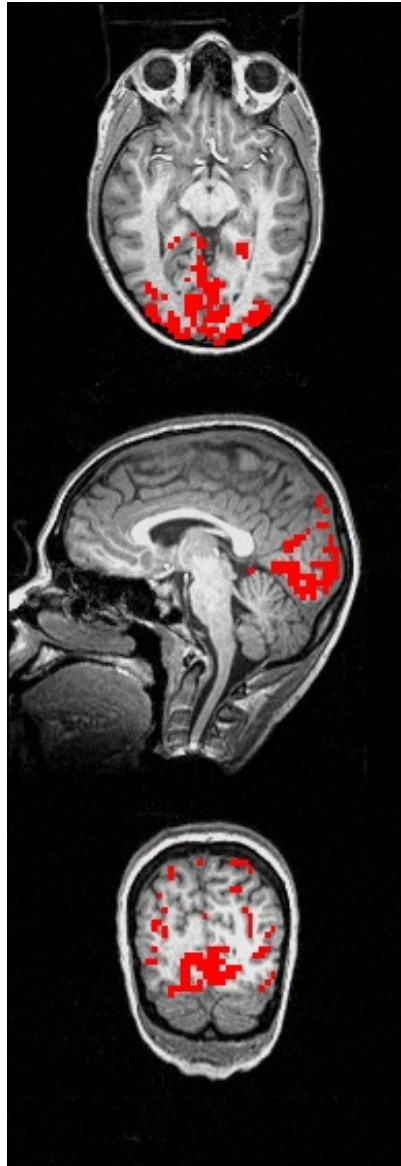


N = 4 error bars: s.e.m.

Decoding Performance



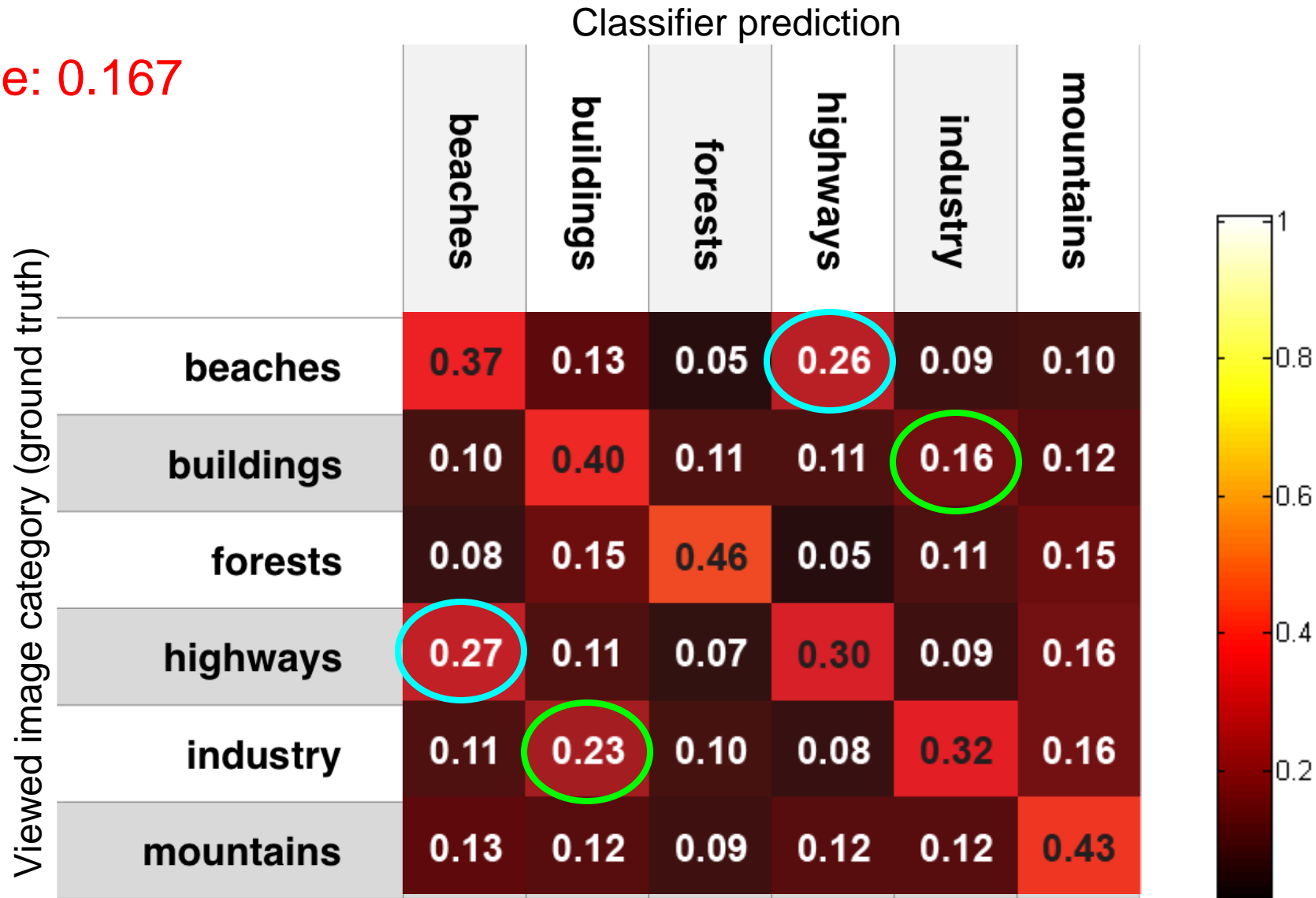
Decoding Performance



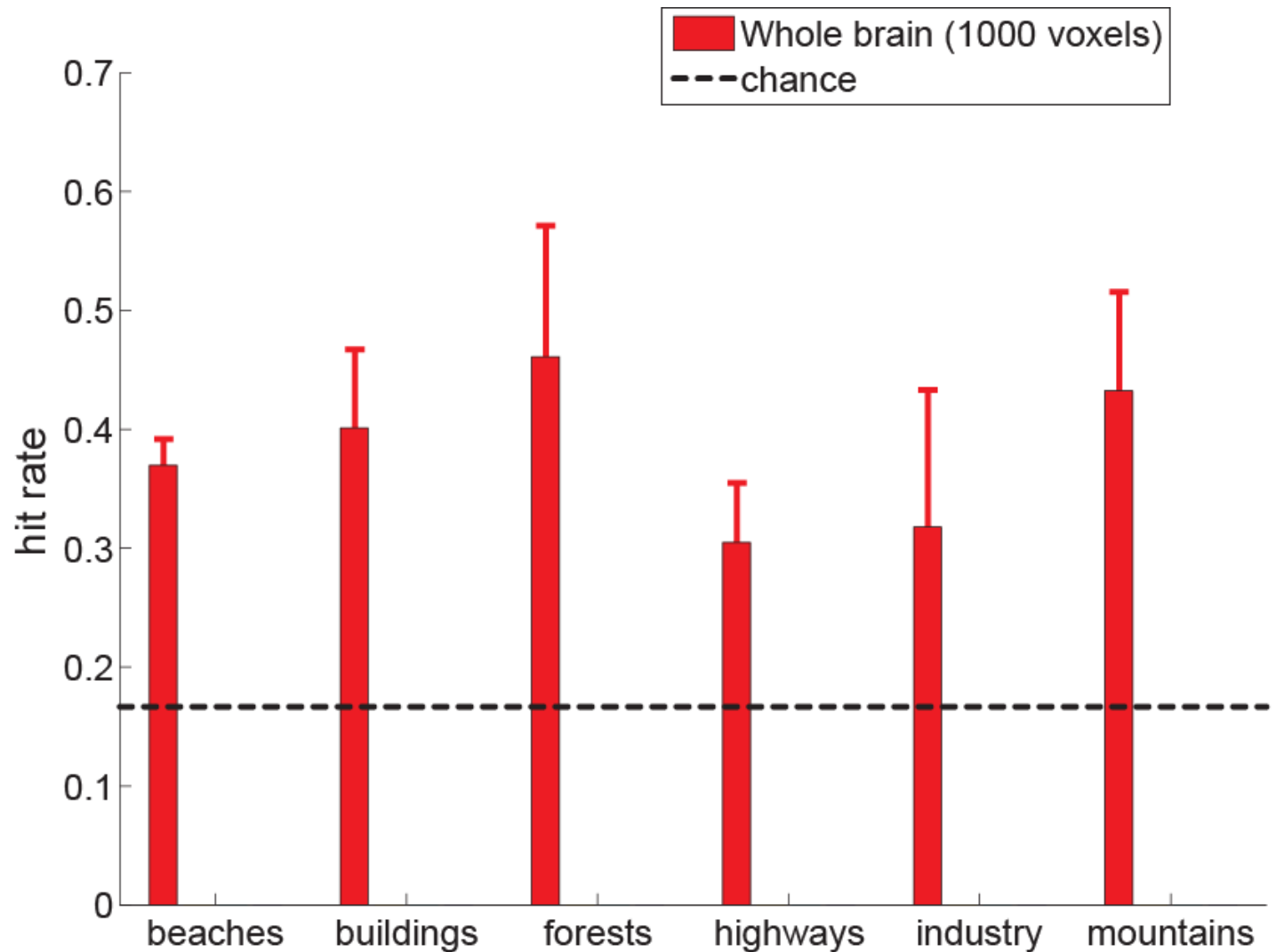
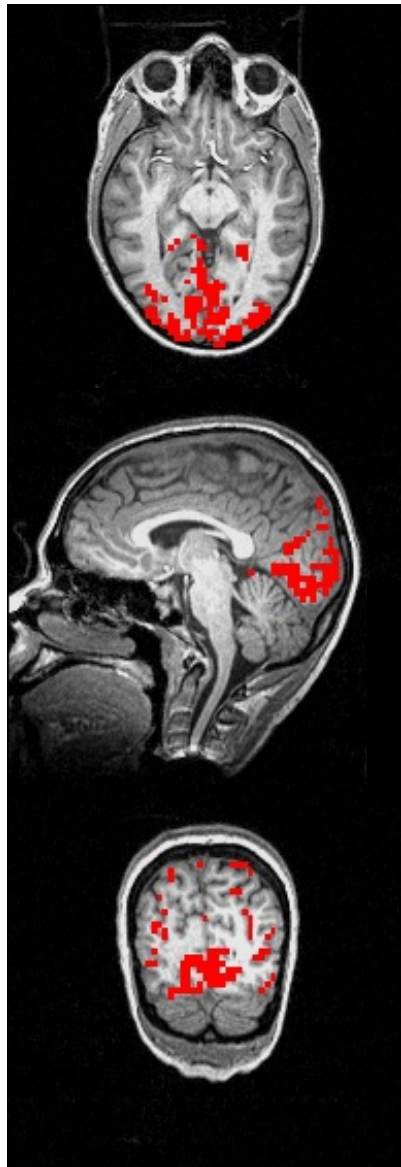
N = 4 error bars: s.e.m.

Decoding Performance

chance: 0.167

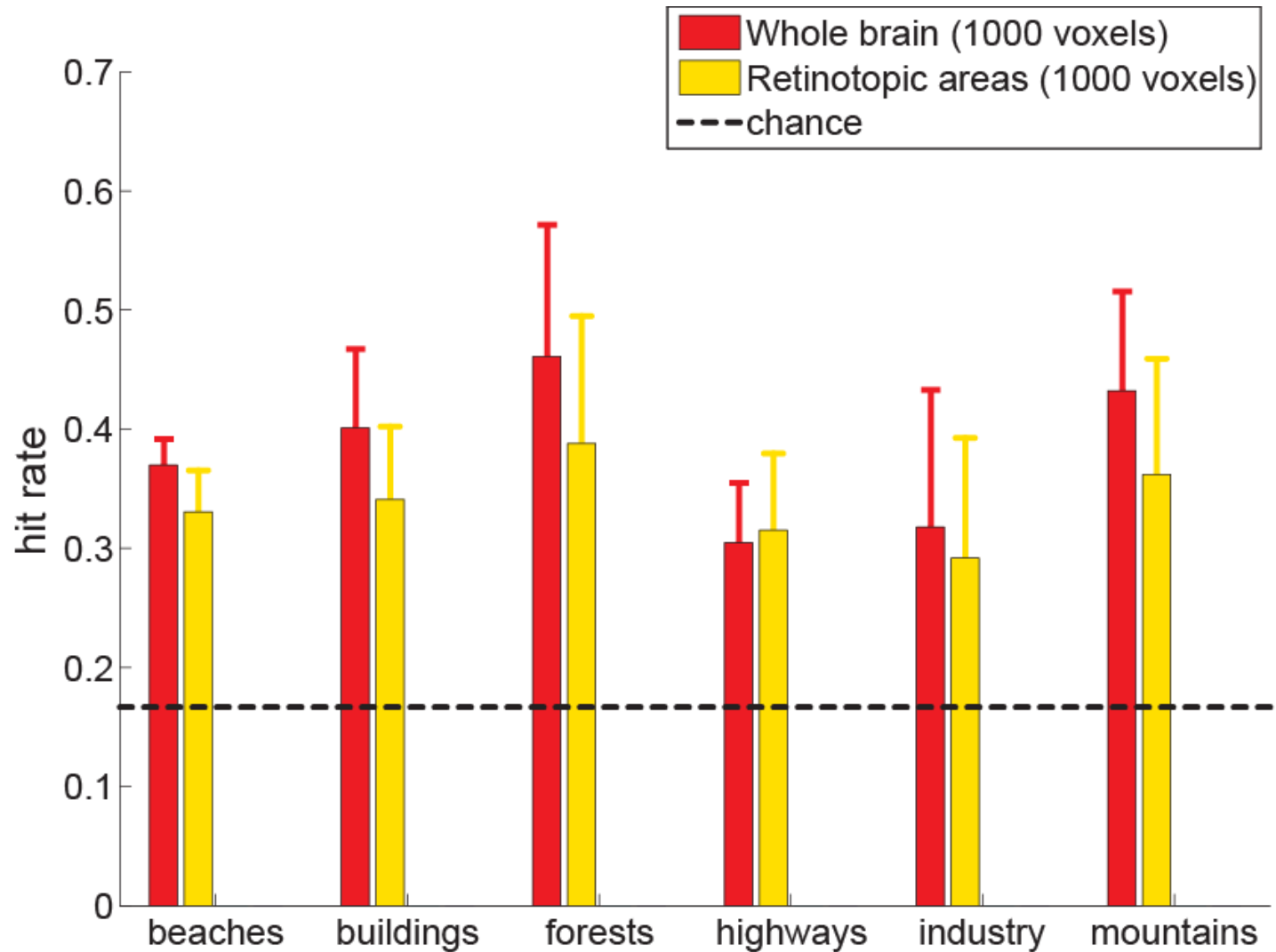
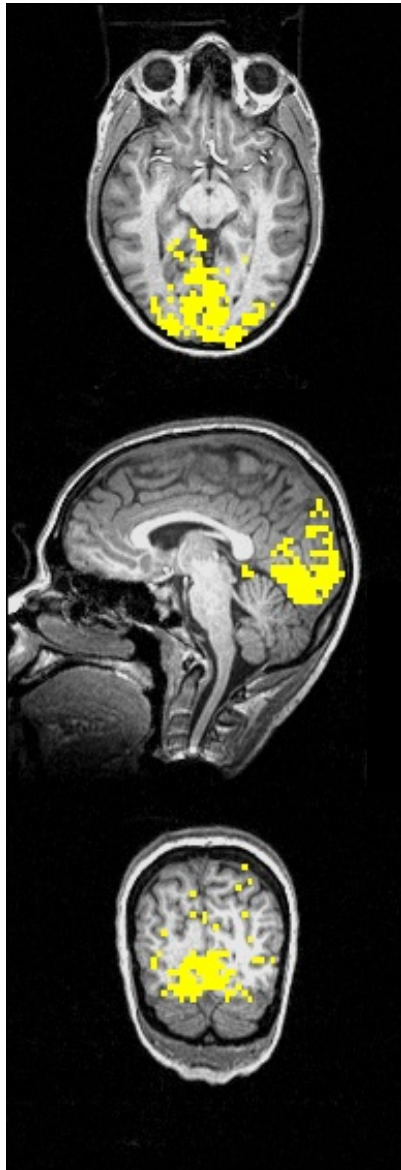


Decoding Performance



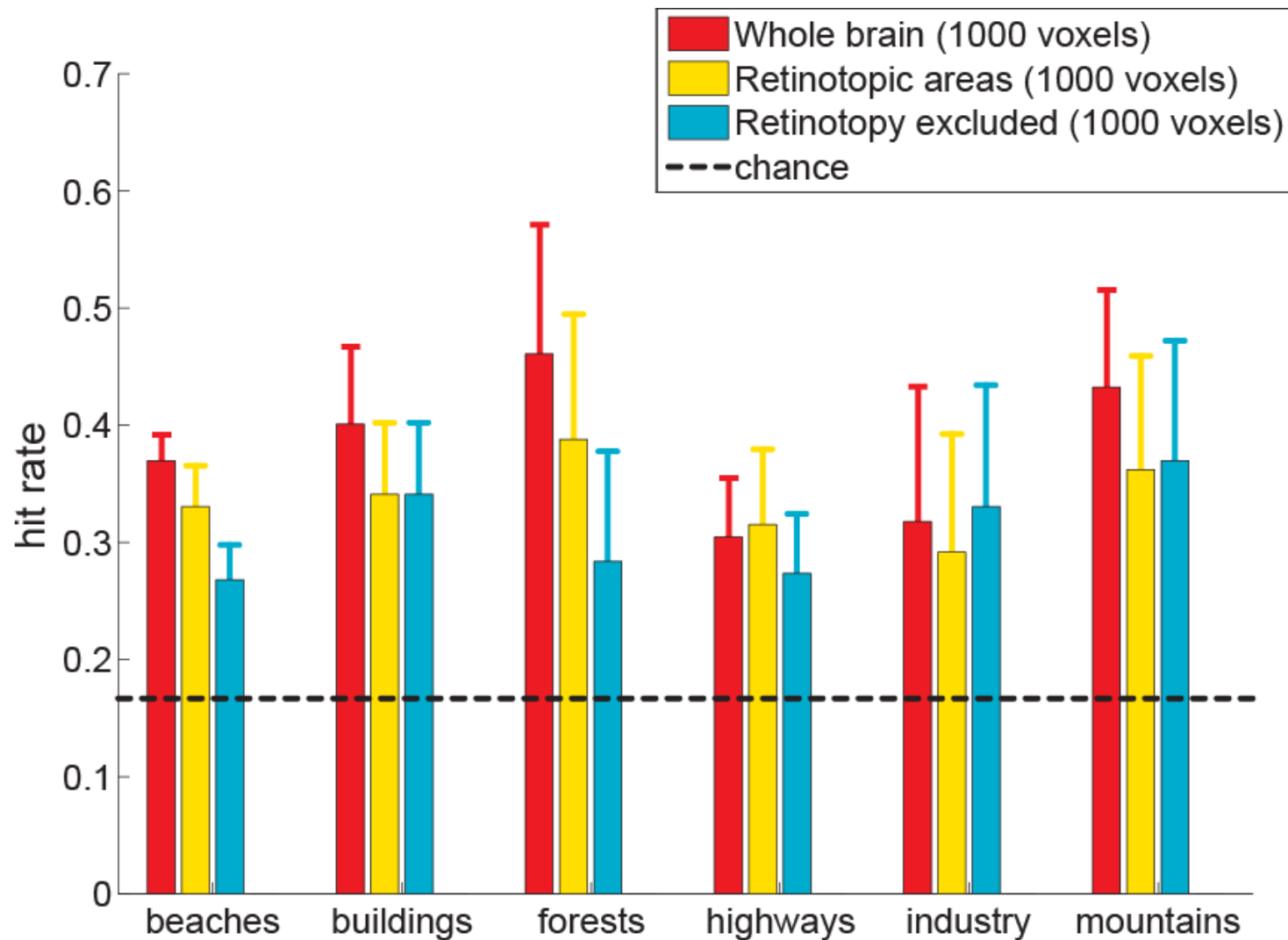
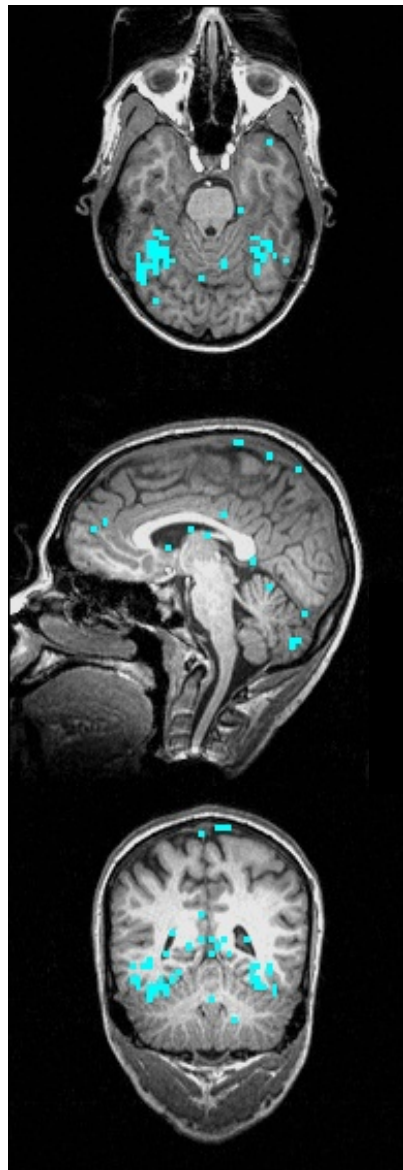
N = 4 error bars: s.e.m.

Retinotopic Areas



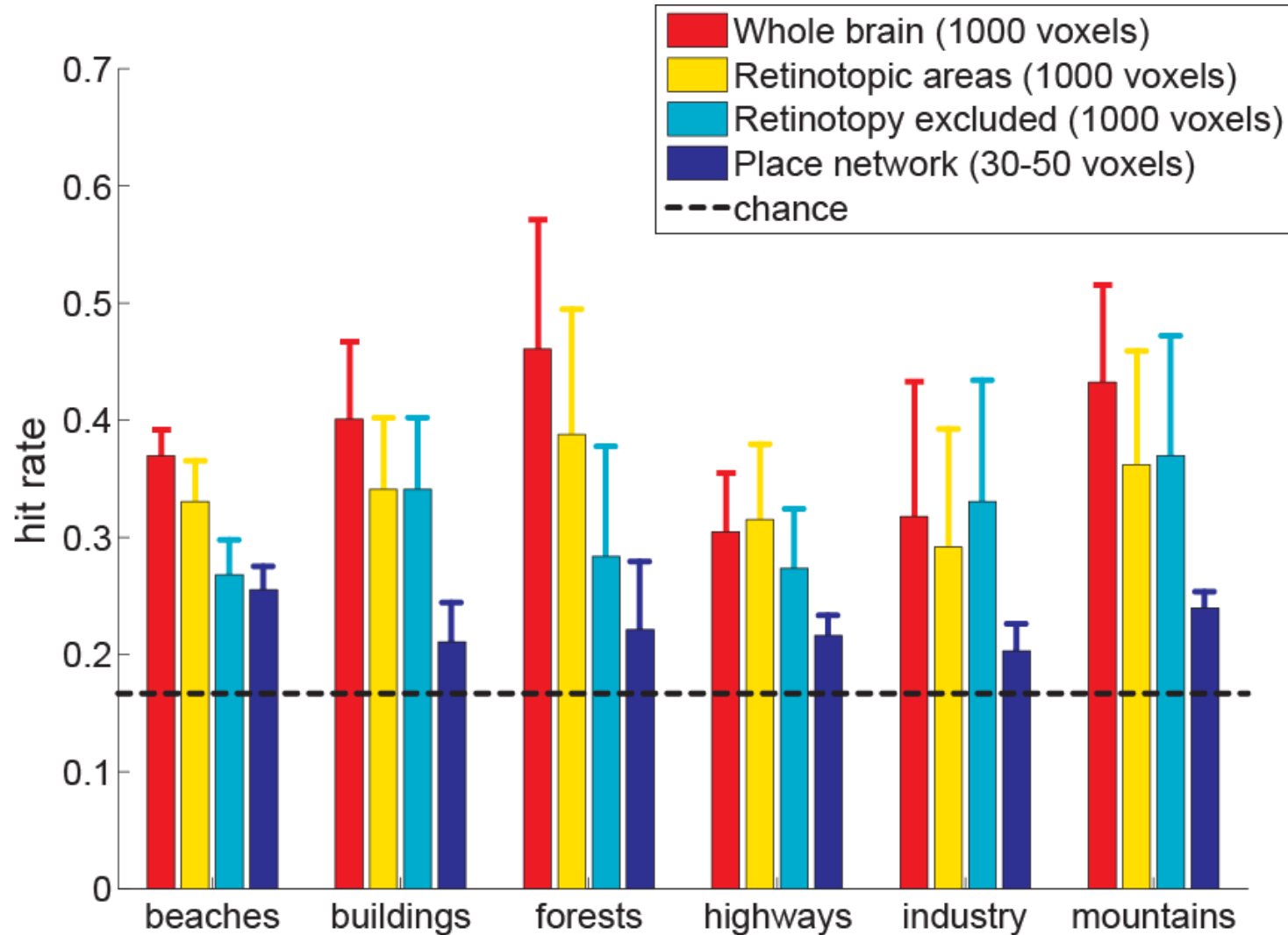
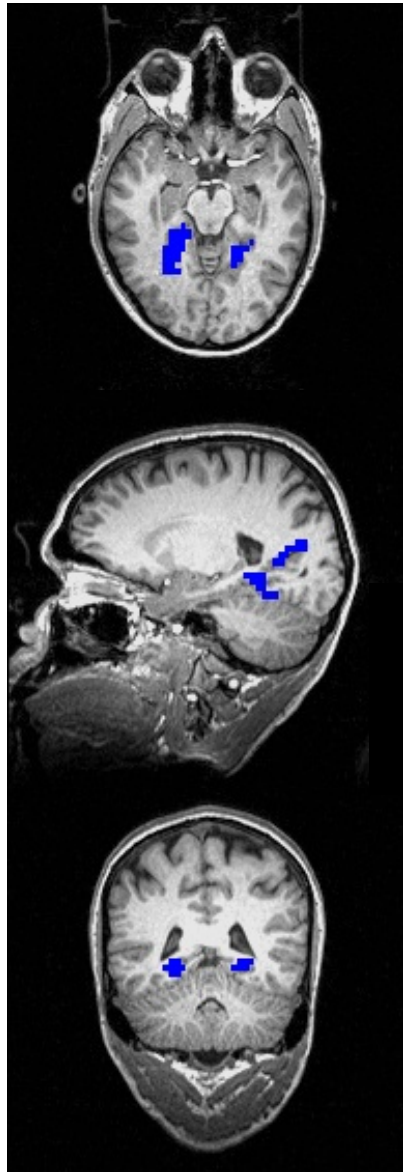
N = 4 error bars: s.e.m.

Retinotopic Areas Excluded

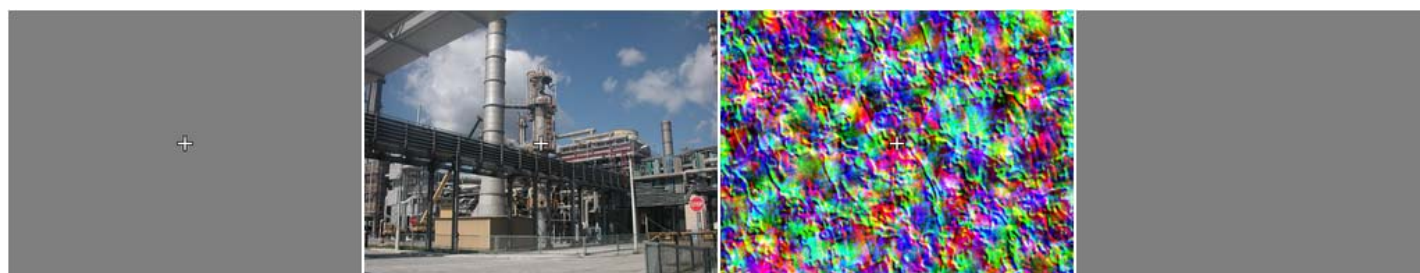


N = 4 error bars: s.e.m.

Place Network (PPA + RSC)



N = 4 error bars: s.e.m.



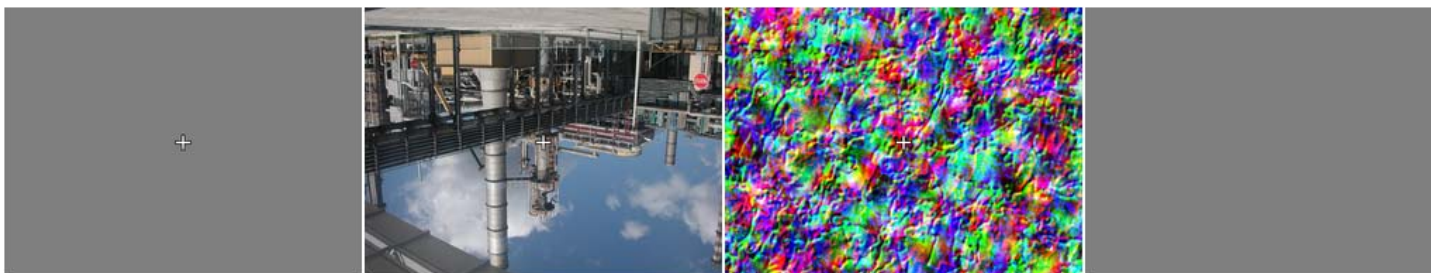
500 ms

32-45 ms

500 ms

< 2000 ms

Upright images



500 ms

32-45 ms

500 ms

< 2000 ms

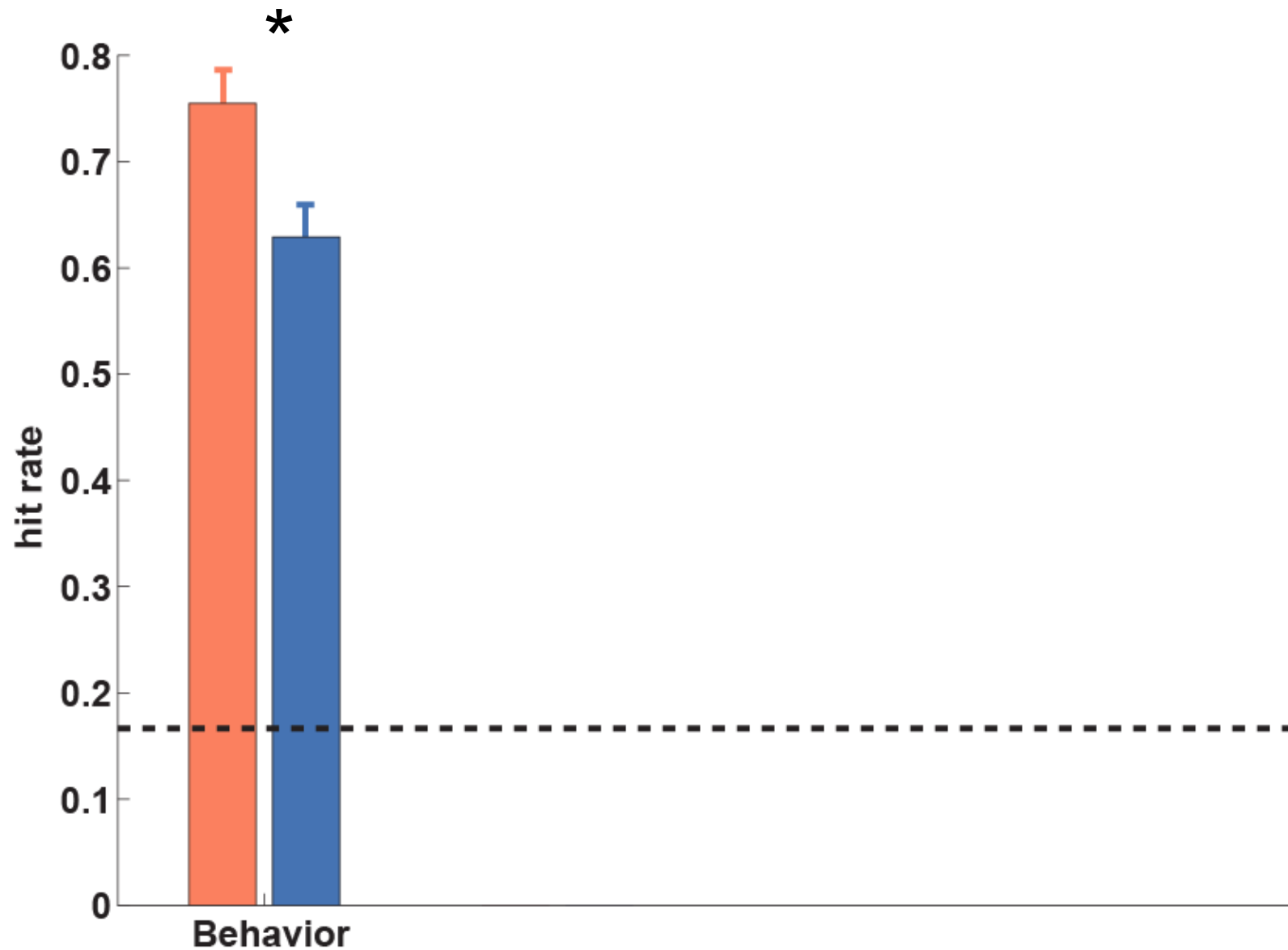
Inverted images

Training: upright only;
Testing: upright & inverted blocks intermixed

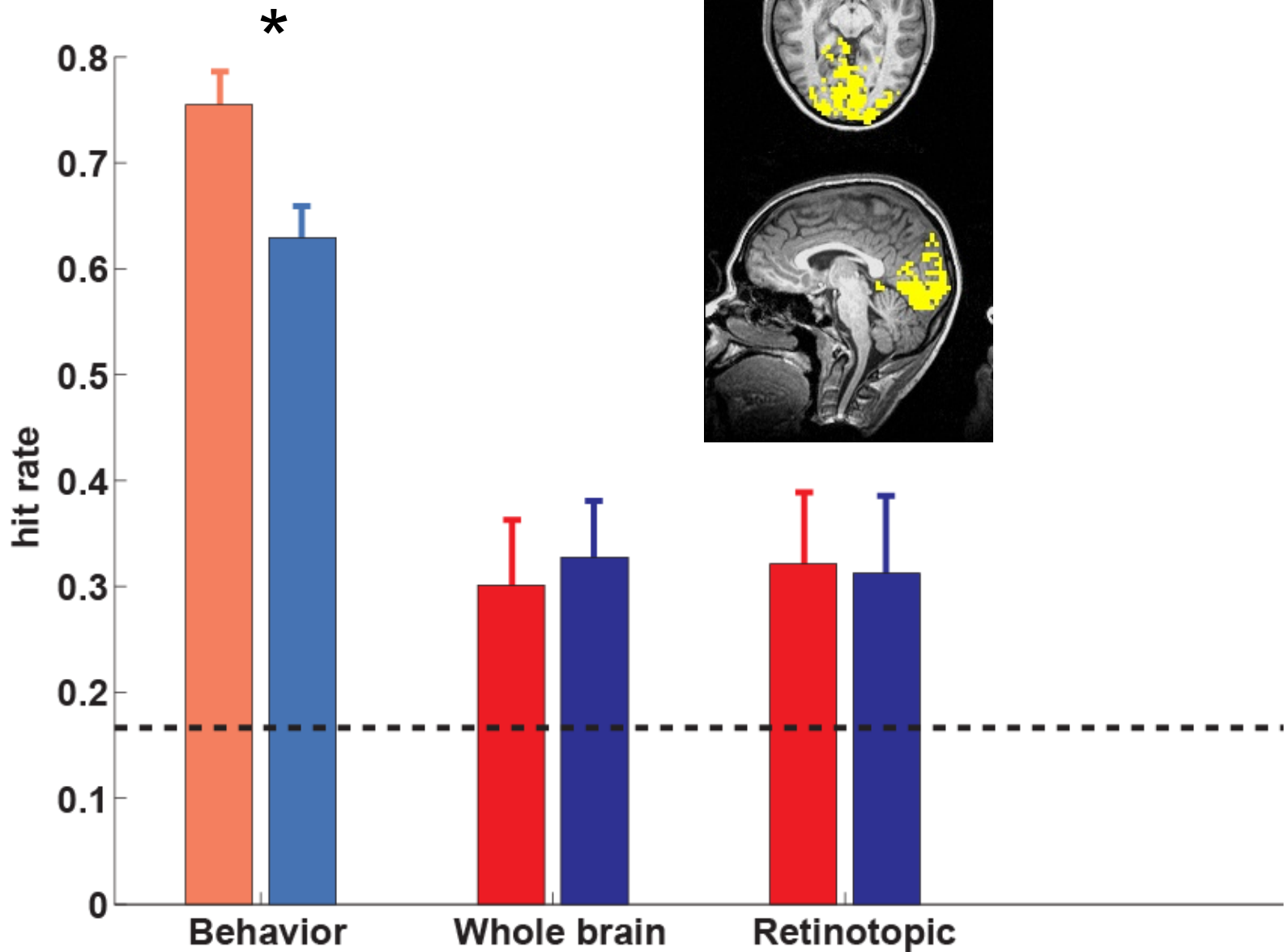


Scene inversion effect

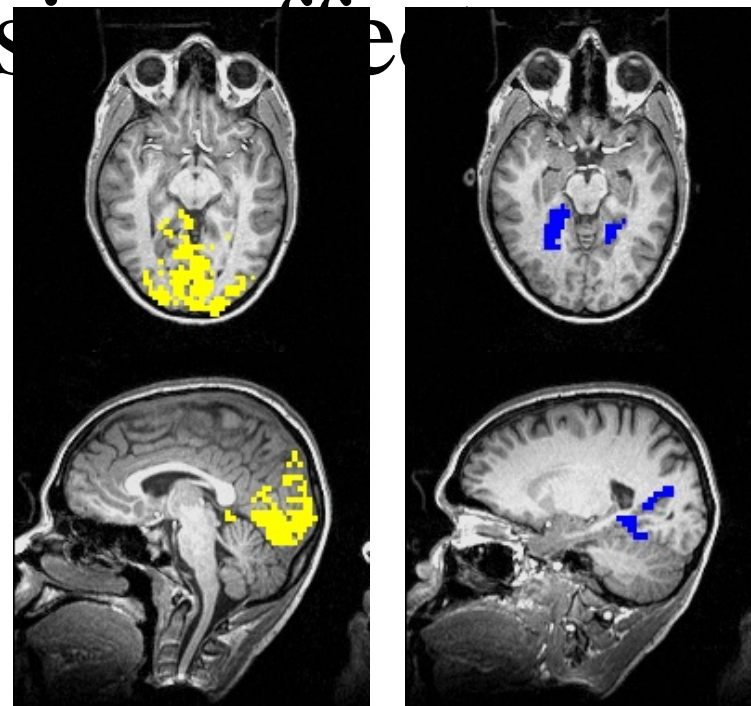
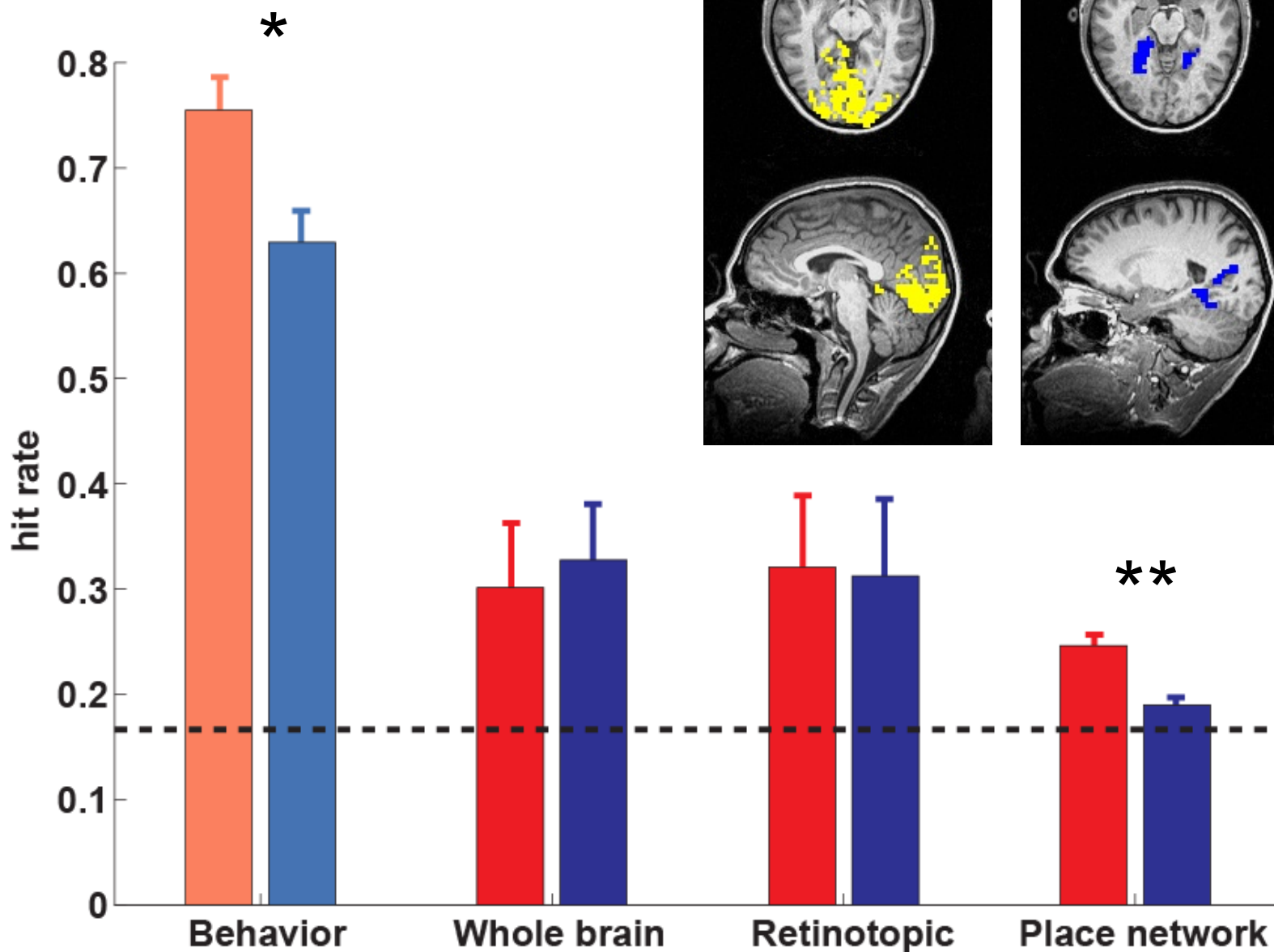
Upright
Inverted
--- chance



Scene inversion effect



Scene inversion effect



Beaches



Highways



Forests



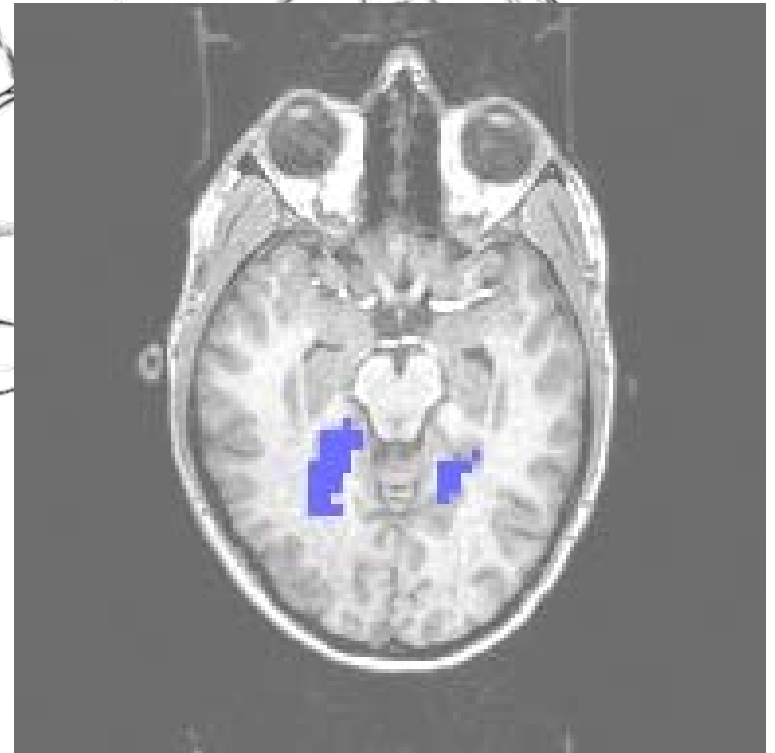
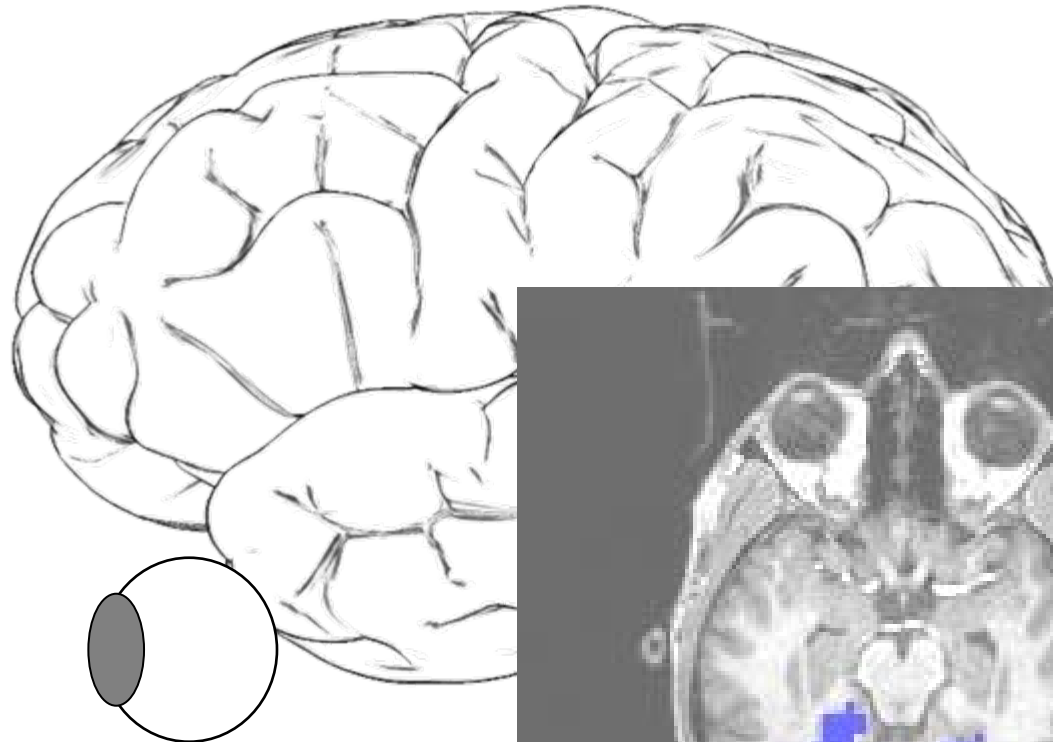
Buildings



Mountains



Industry



Beaches



Highways



Forests



Buildings



Mountains



Industry



livingroom bedroom



highway



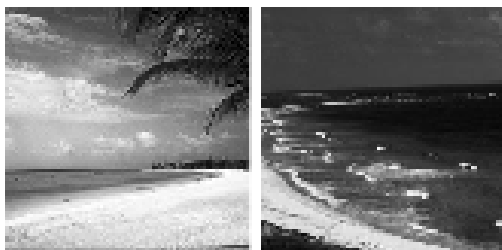
tall bldg



suburb



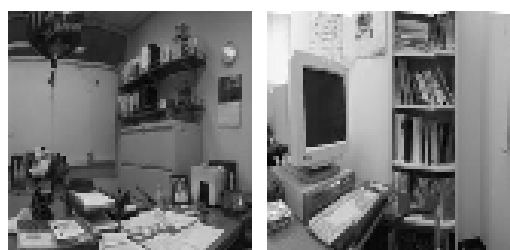
coast



kitchen



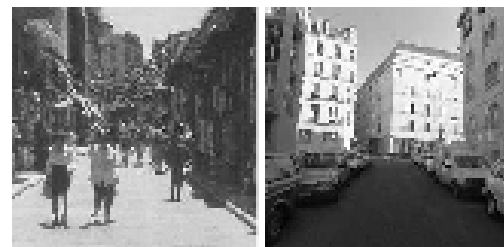
office



ins. city



streets



forest



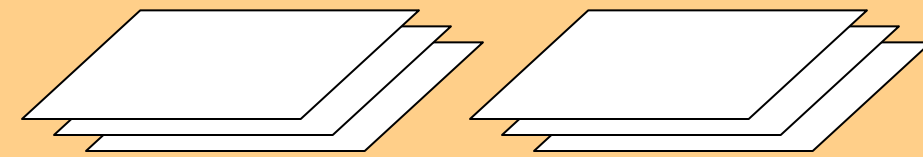
mountain



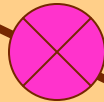
o. country



learning



feature detection
& representation



codewords dictionary

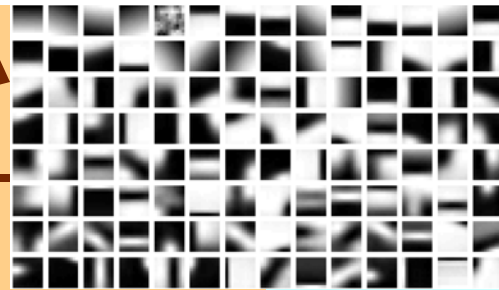
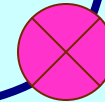
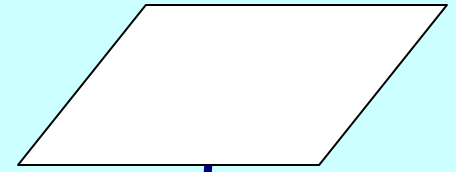


image representation



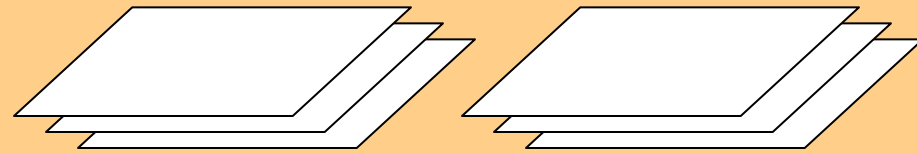
**category models
(and/or) classifiers**

recognition

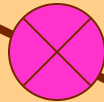


**category
decision**

Representation



1. feature detection
& representation



2. **codewords dictionary**

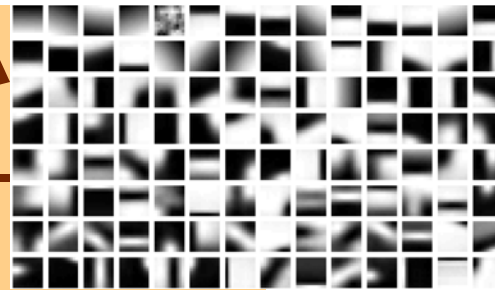
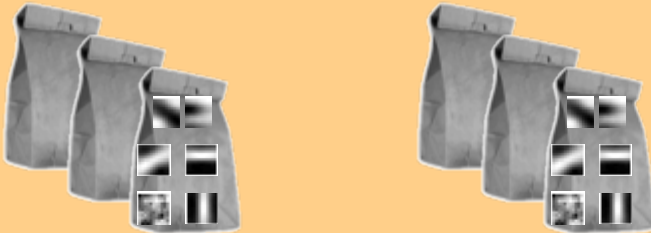
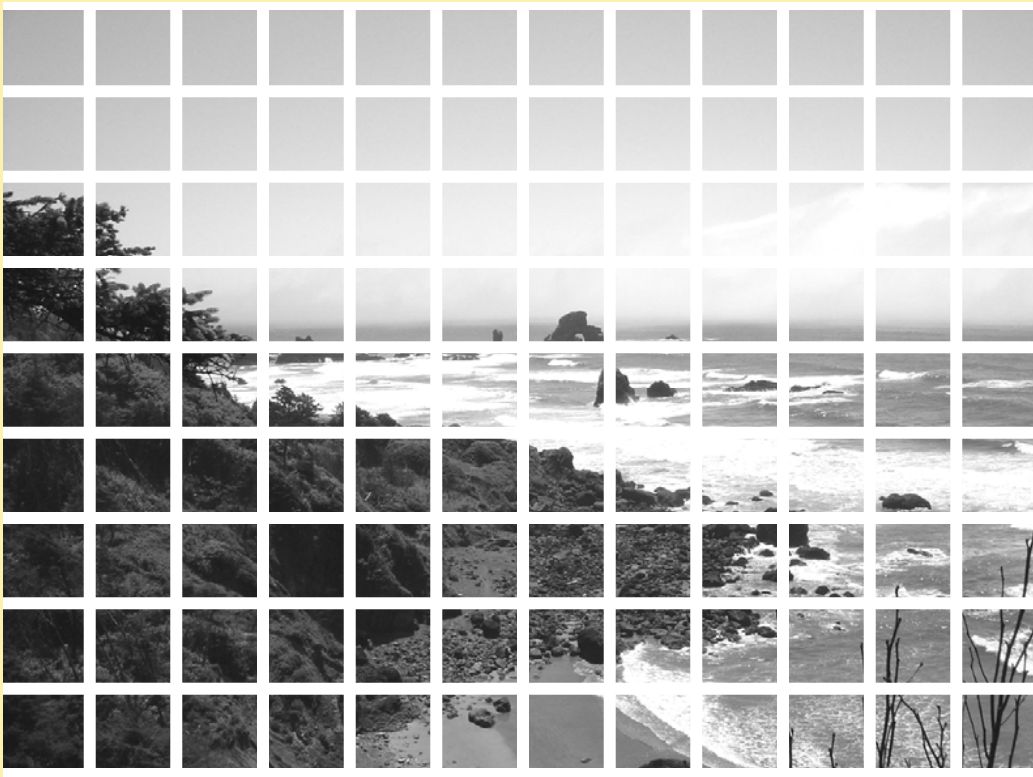


image representation

3.



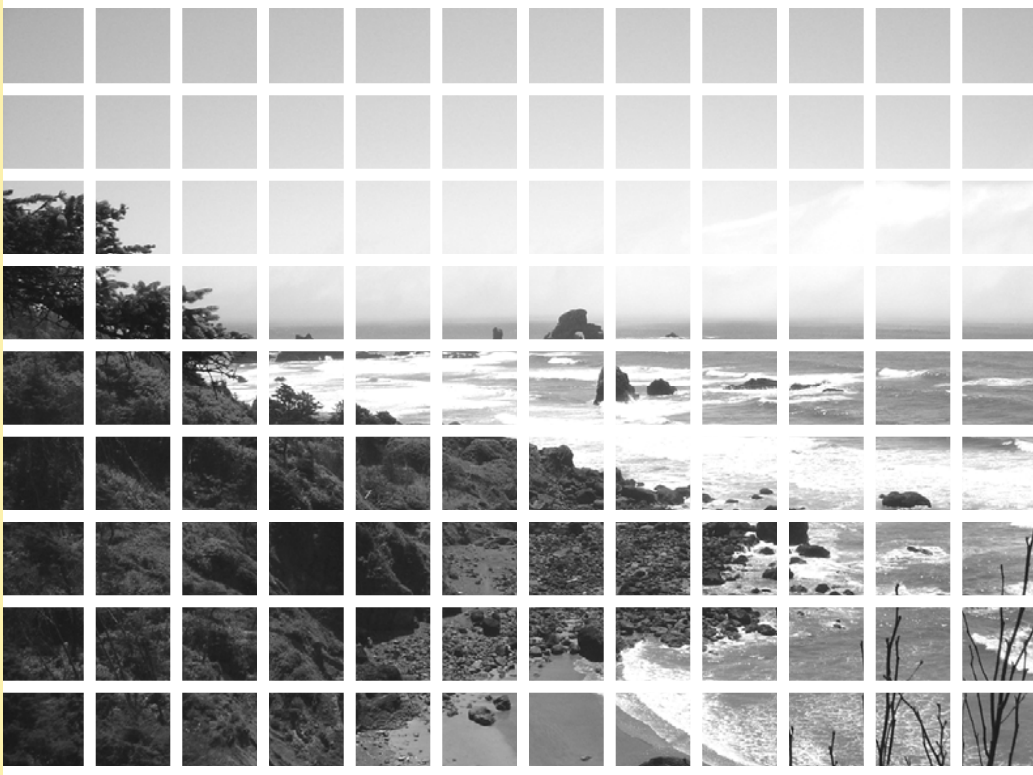
1.Feature detection and representation



**extract
interest points**

- DoG
- Saliency detector (Kadir and Brady)
- grid

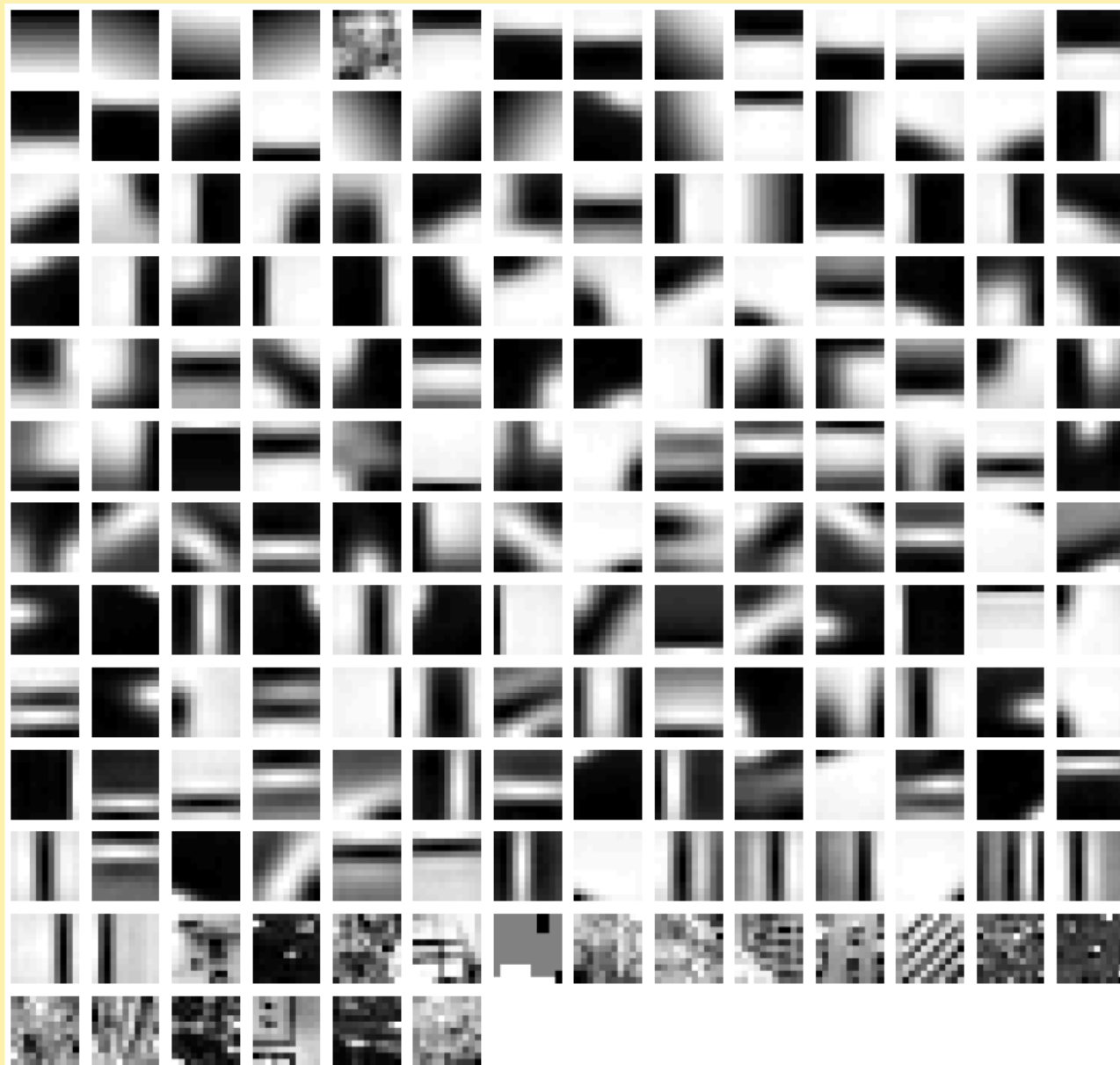
1.Feature detection and representation



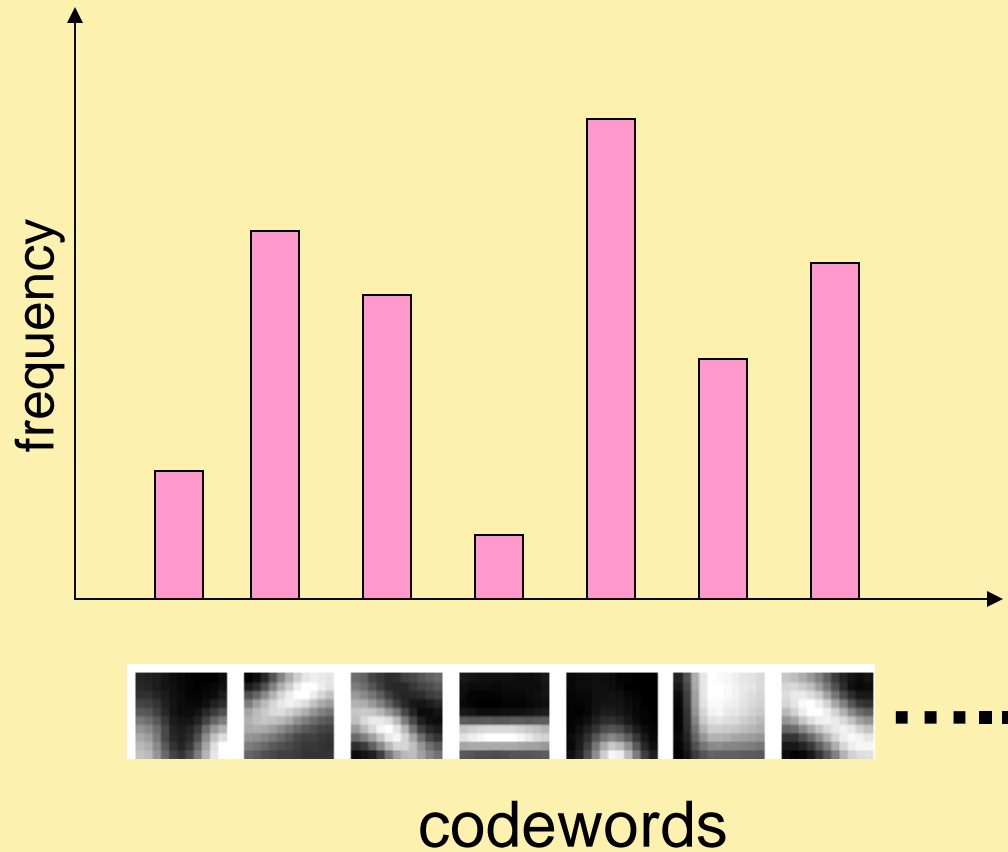
**represent
interest points**

- SIFT (Lowe '99)
- gray scale values

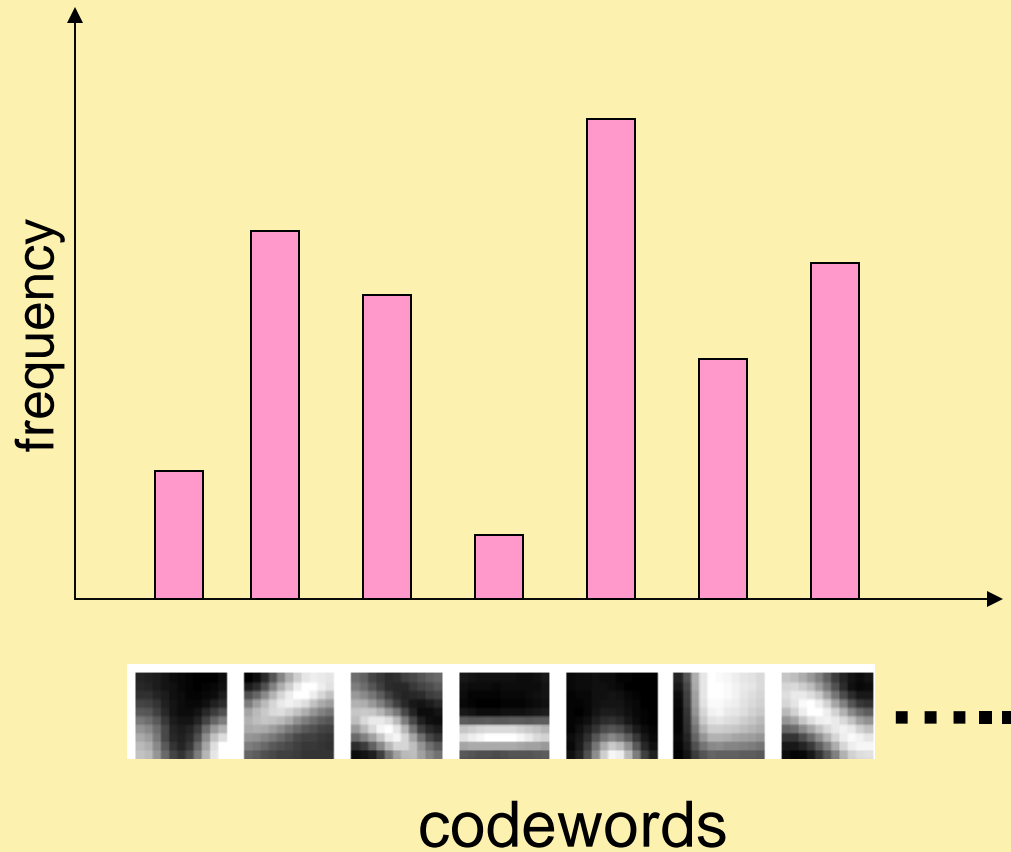
2. Codewords dictionary formation



3. Image representation



3. Image representation



Analogy to documents

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our eyes.

For a long time, the retinal image was considered as a movie screen. It is now known that the image is processed in a more complex way. Following the discovery of the pathway to the various centers of the cortex, Hubel and Wiesel have demonstrated that the message about the image falling on the retina undergoes a step-by-step analysis in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.

**sensory, brain,
visual, perception,
retinal, cerebral cortex,
eye, cell, optical
nerve, image
Hubel, Wiesel**

China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports to \$750bn, compared with \$580bn in 2004, and a fall in imports to \$660bn. The ministry said the surplus would annoy the US.

China's government has deliberately agreed to a trade surplus with the US, the yuan is valued at 8 yuan to the dollar, the government also needs to keep the yuan's value low to stimulate demand so that it can keep the economy growing. China has been allowed to trade the yuan against the dollar since 2005 and permitted it to trade within a narrow band but the US wants the yuan to be allowed to trade freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

**China, trade,
surplus, commerce,
exports, imports, US,
yuan, bank, domestic,
foreign, increase,
trade, value**

learning



feature detection
& representation



codewords dictionary

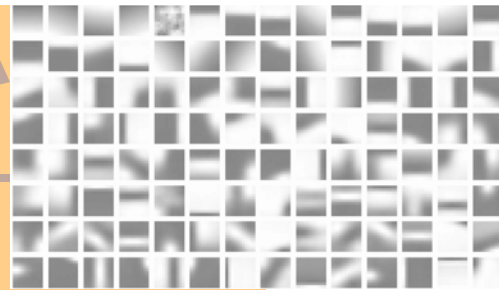
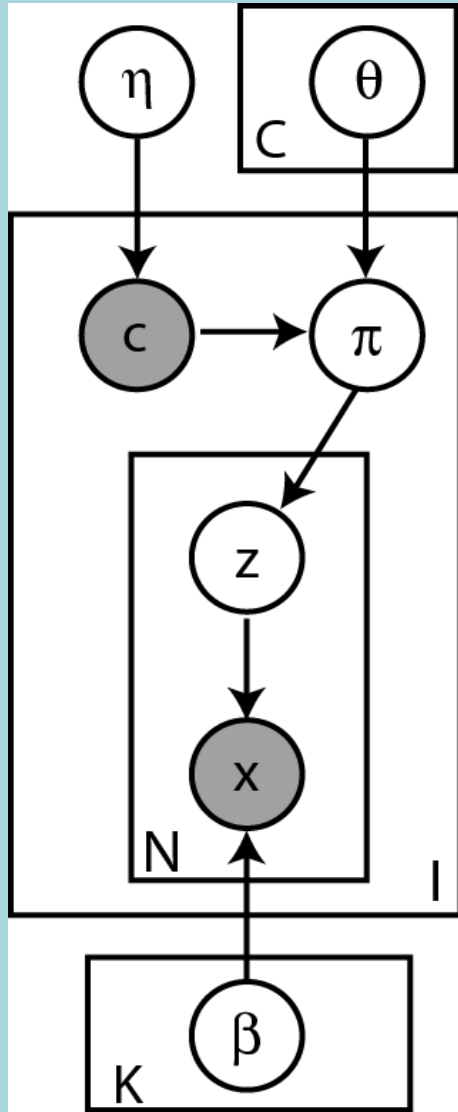


image representation

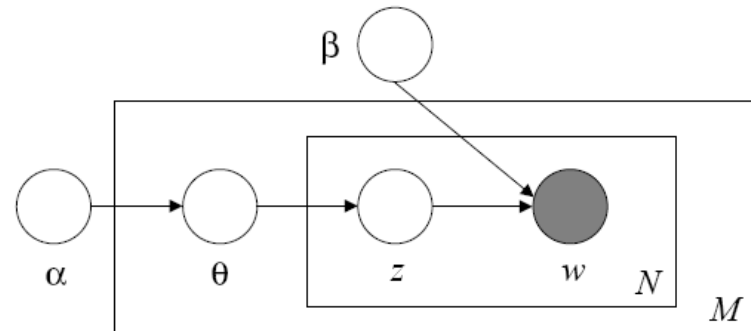


**category models
(and/or) classifiers**

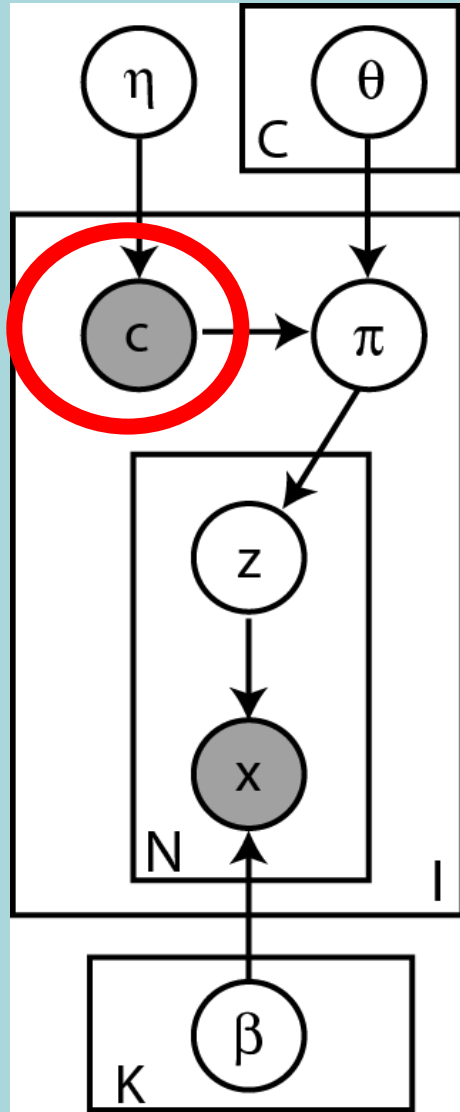
A Generative Model



LDA: Blei, Ng, & Jordan. 2003



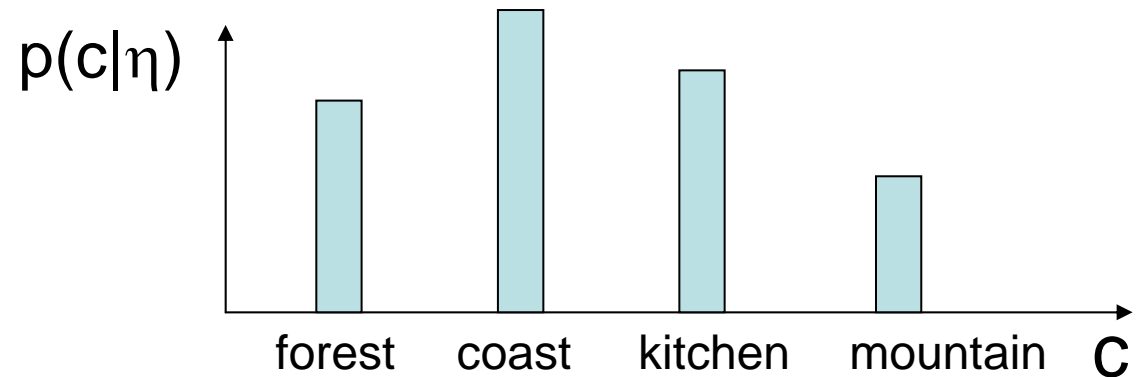
A Generative Model



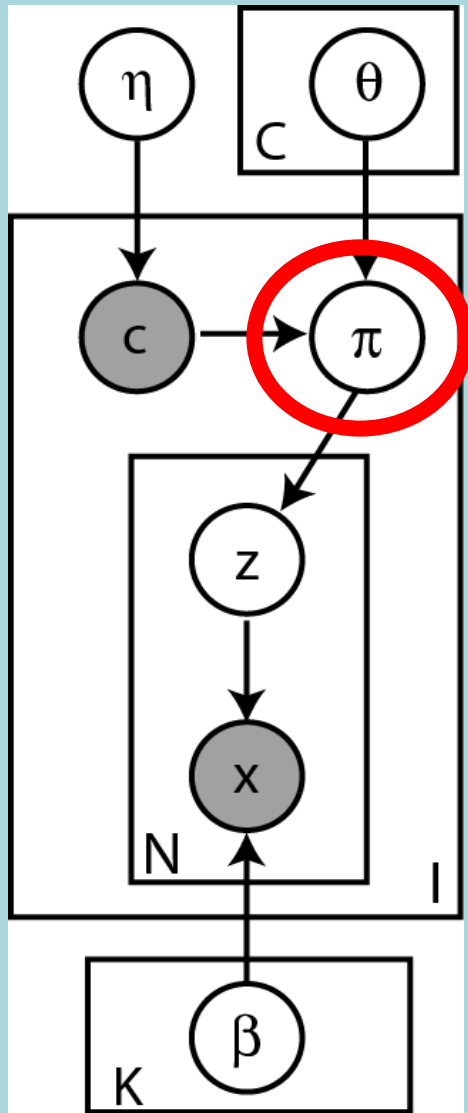
scene category



discrete variable: $c \sim p(c|\eta)$



A Generative Model



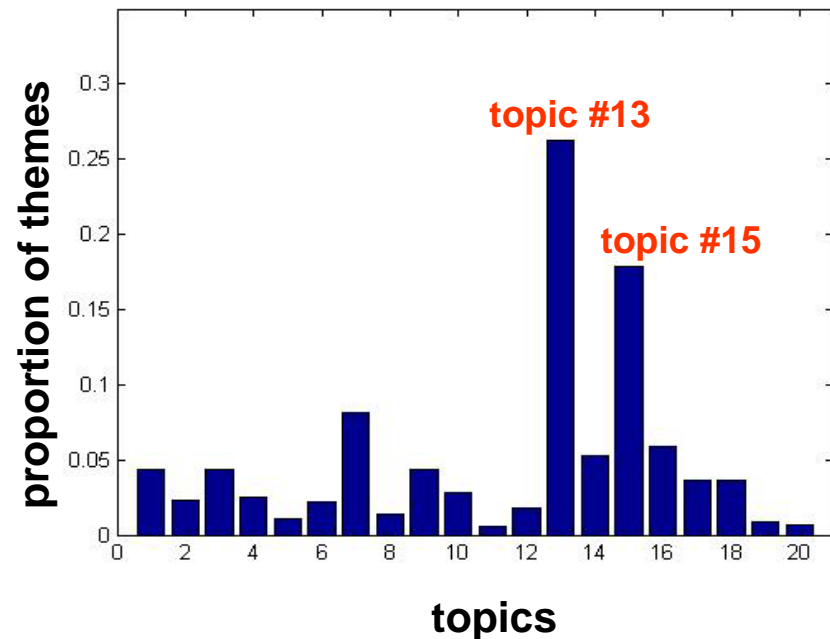
mixing parameter for the latent topics



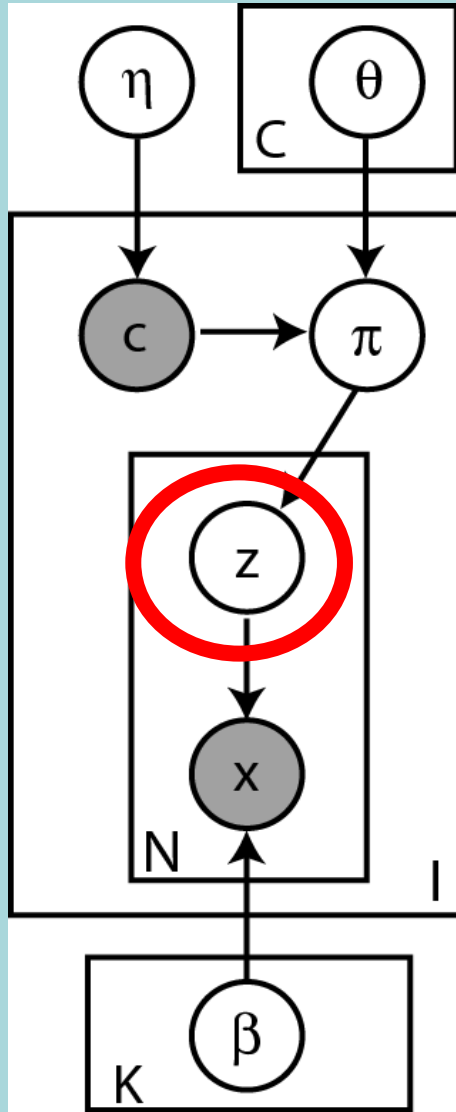
$$\pi \sim p(\pi | c, \theta)$$

$$\sim \text{Dir}(\pi | c, \theta)$$

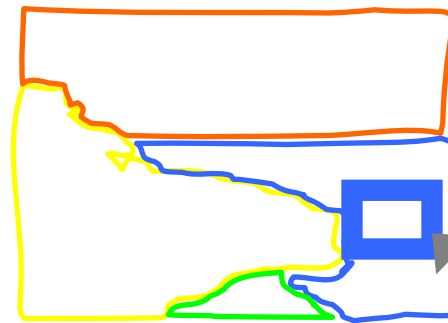
$$\text{where } \sum_{k=1}^K \pi_k = 1$$



A Generative Model



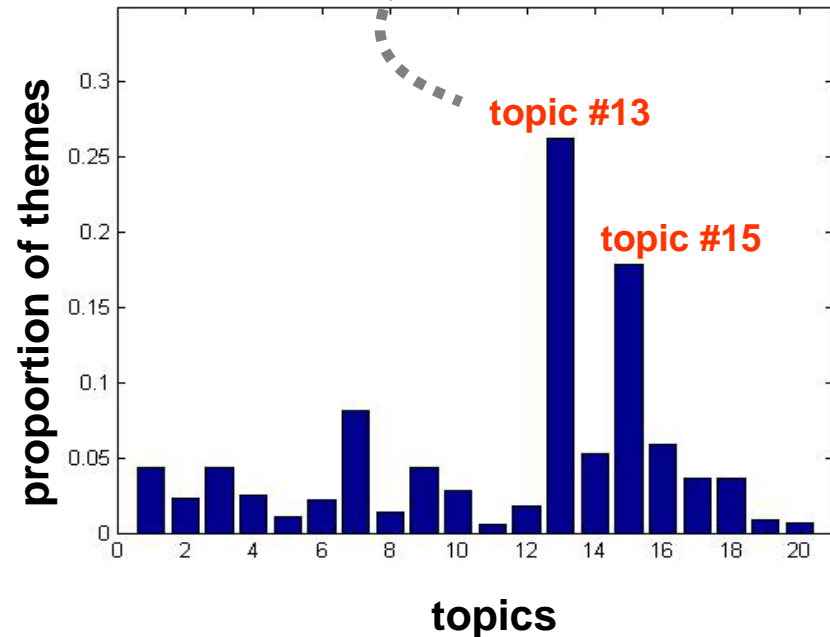
topic label



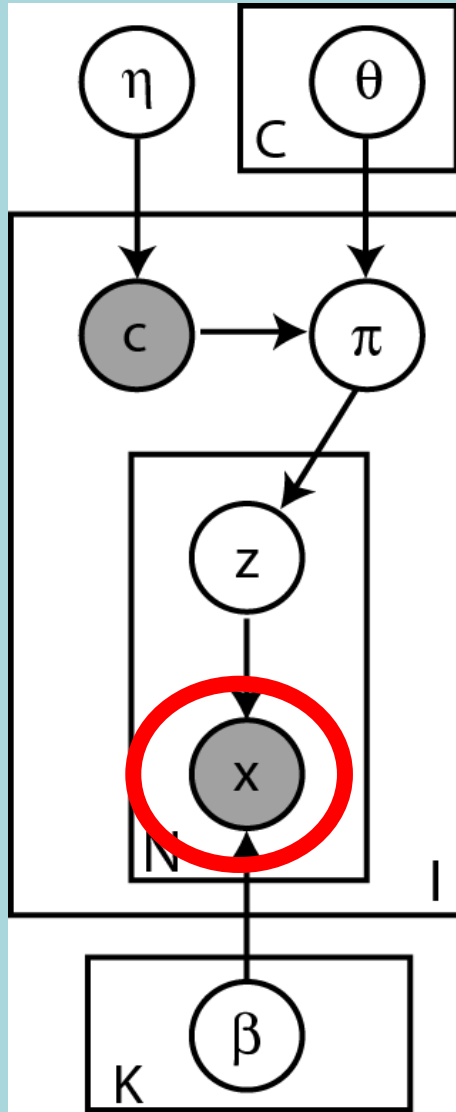
discrete variable:

$$z \sim p(z|\pi)$$

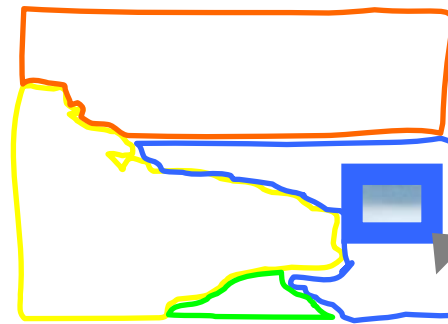
$$\sim \text{Mult}(z|\pi)$$



A Generative Model



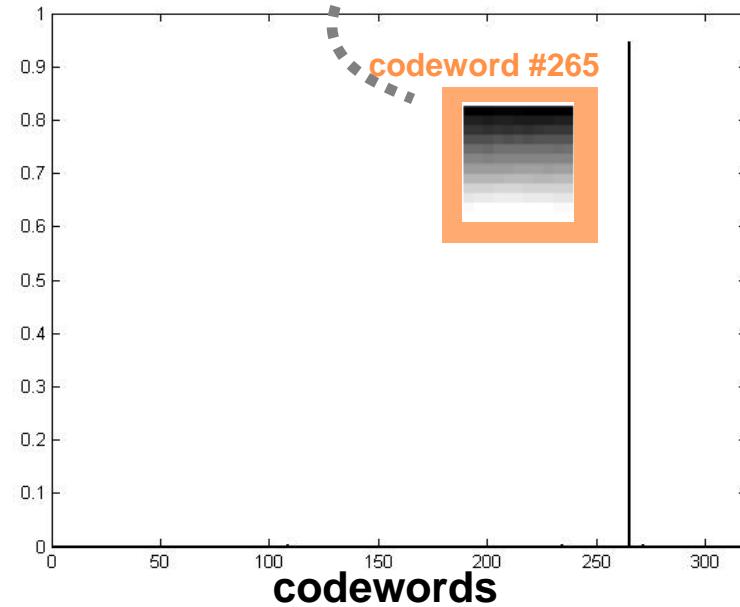
patch label



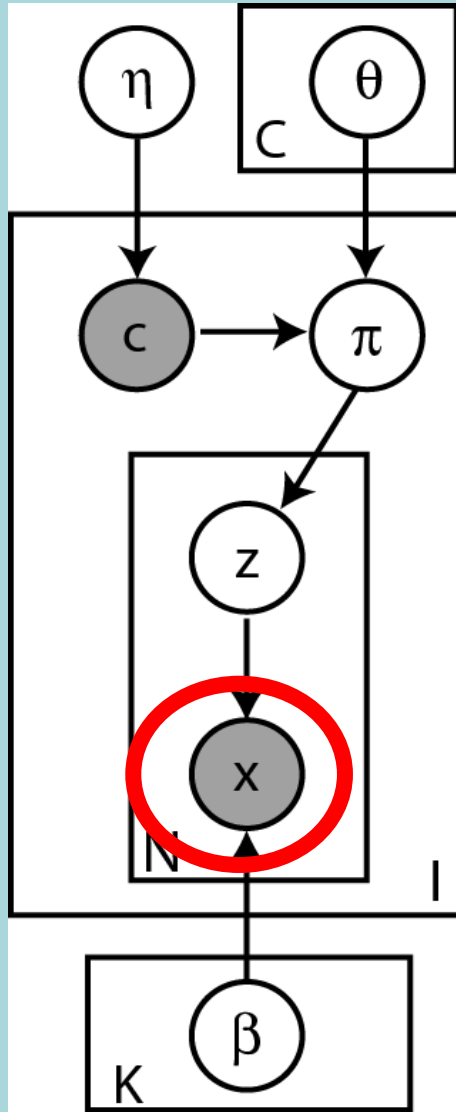
discrete variable:

$$x \sim p(x|z, \beta)$$
$$\sim \text{Mult}(x|z, \beta)$$

expected value of β given 'z=13'

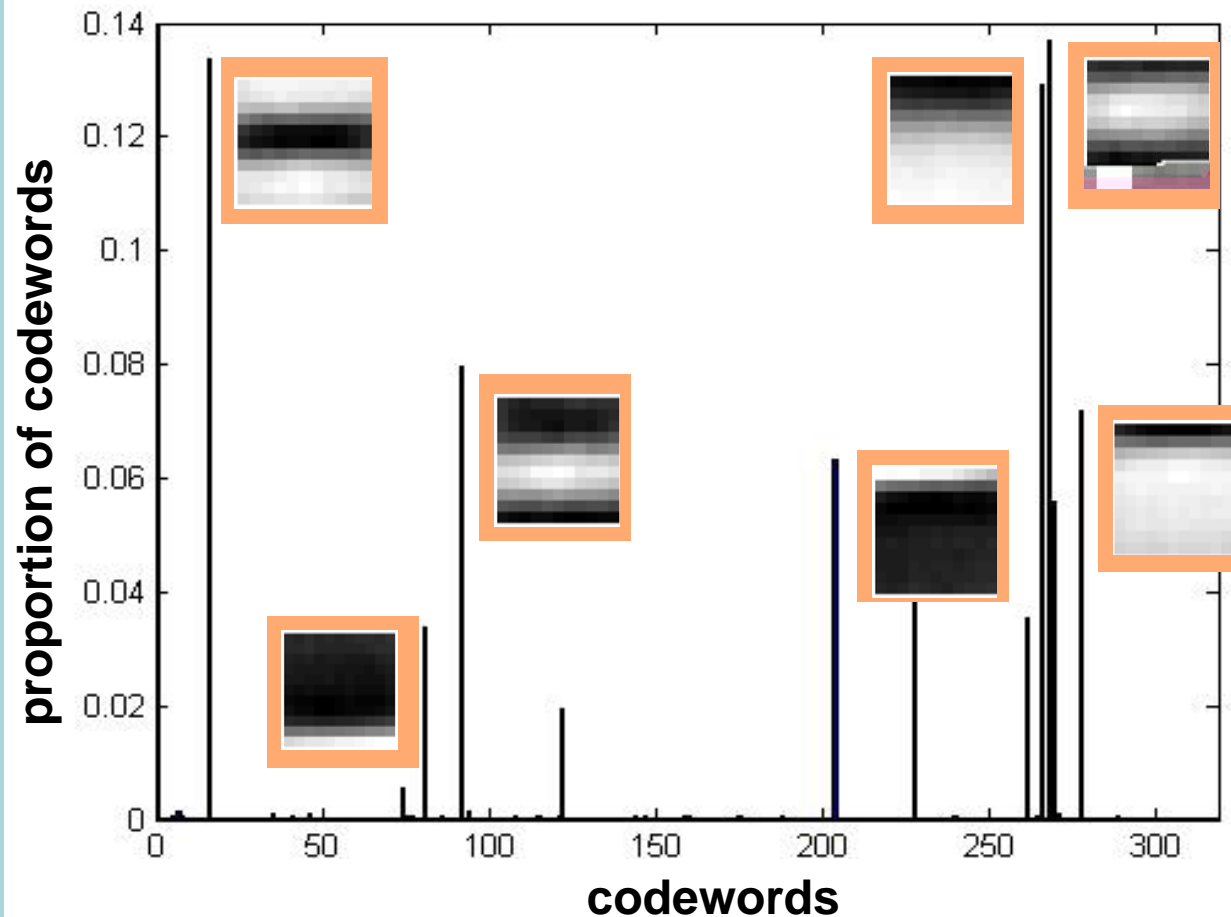


A Generative Model

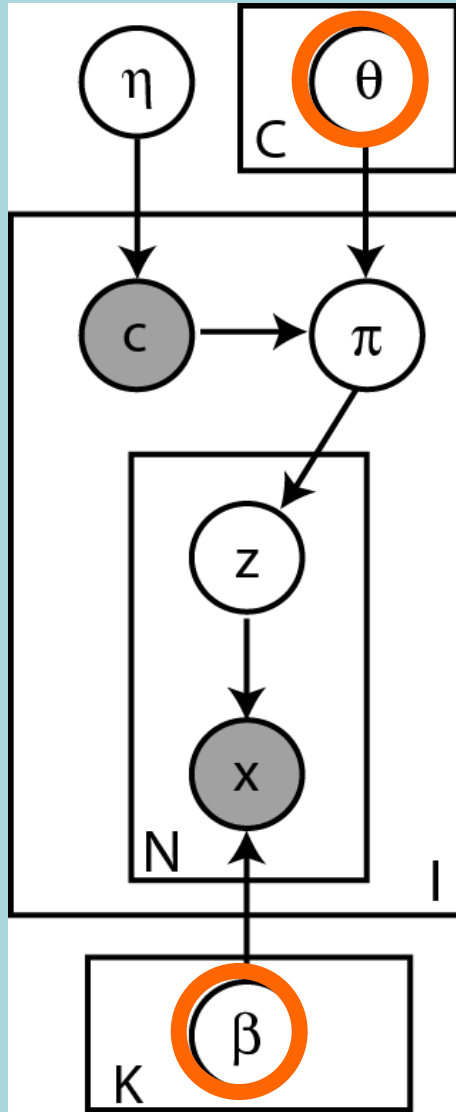


patch label

expected value of β given 'z=15'



A Generative Model



learning

Find the 'best' θ and β

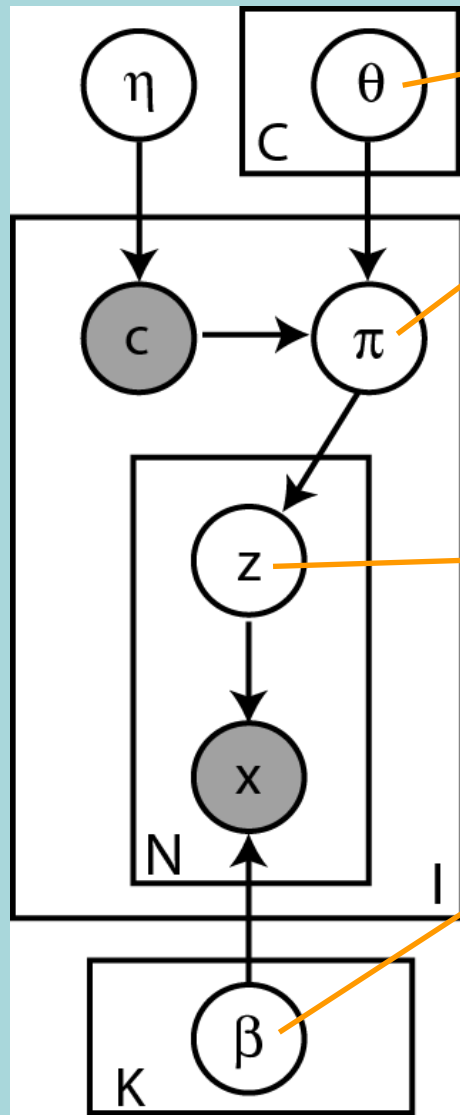
joint probability

$$p(x, z, \pi | \theta, \beta, c) = p(\pi | c, \theta) \prod_n^N p(z_n | \pi) p(x_n | z_n, \beta)$$

$$p(x | \theta, \beta, c) = \int p(\pi | c, \theta) \left(\prod_n^N \sum_{z_n} p(z_n | \pi) p(x_n | z_n, \beta) \right) d\pi$$

- exact inference is intractable
- use Variational Inference

A Generative Model



Variational Inference

Maximum Likelihood estimation (Minka 2000)

$$\gamma_{ck} = \theta_{ck}^0 + \sum_n \langle \delta(z_n^k = 1) \rangle$$

$$\langle \log \pi_{ck} \rangle = \Psi(\gamma_{ck}) - \Psi\left(\sum_k \gamma_{ck}\right)$$

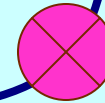
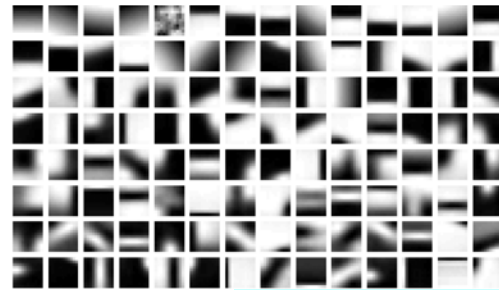
$$\langle \delta(z_n^k = 1) \rangle = \exp\left\{ \langle \log \pi_{ck} \rangle + \sum_t \langle \log \beta_{kt} \rangle \delta(x_n^t = 1) \right\}$$

$$\xi_{kt} = \zeta^0 + \sum_i \sum_n \langle \delta(z_{i,n}^k = 1) \rangle \delta(x_{i,n}^t = 1)$$

$$\langle \log \beta_{kt} \rangle = \Psi(\xi_{kt}) - \Psi\left(\sum_t \xi_{kt}\right)$$

Recognition

codewords dictionary

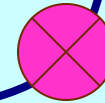
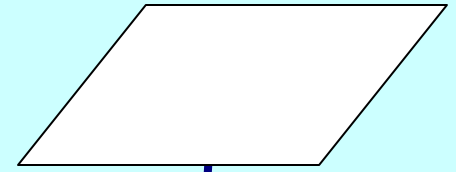
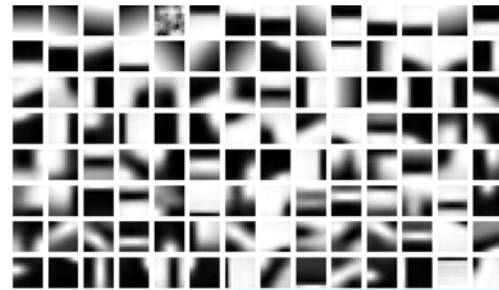


**category models
(and/or) classifiers**

**category
decision**

Recognition

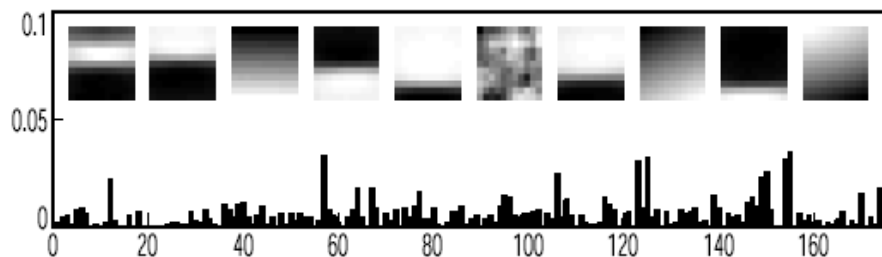
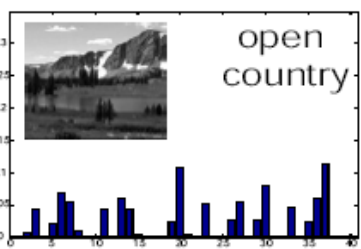
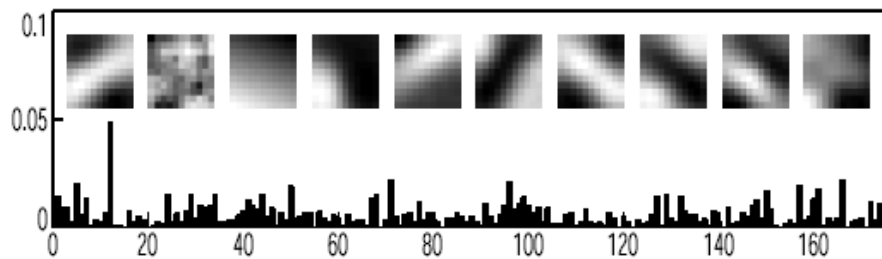
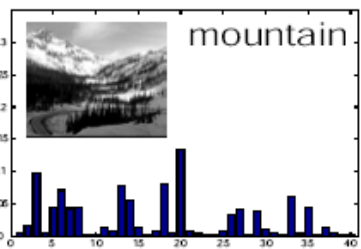
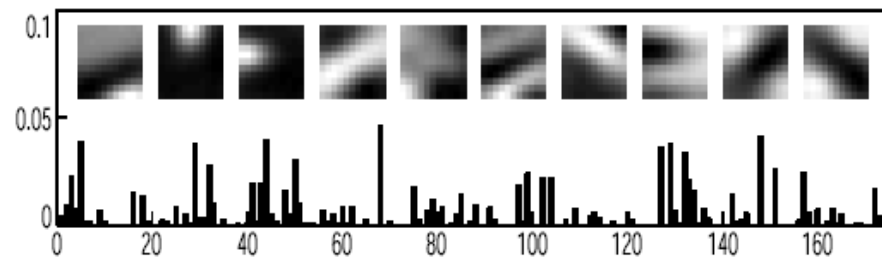
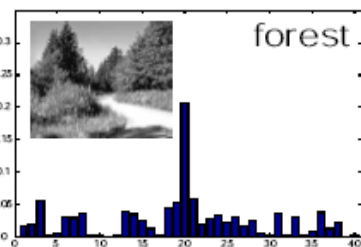
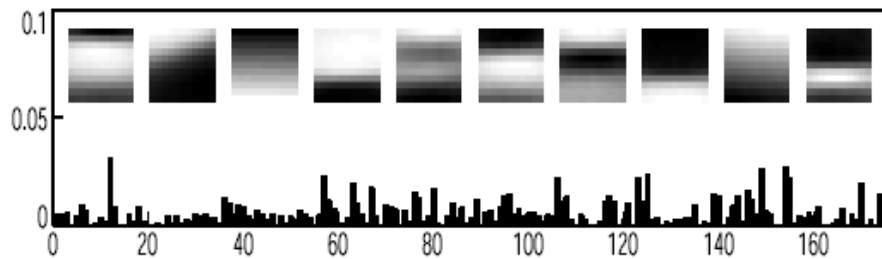
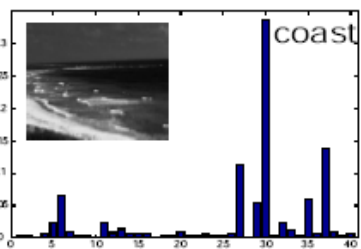
codewords dictionary

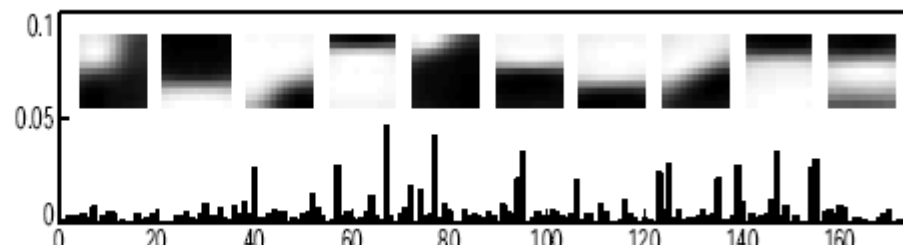
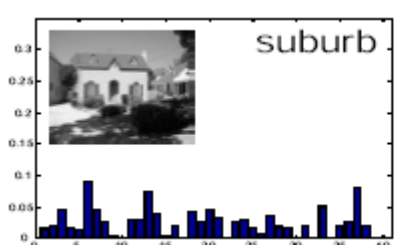
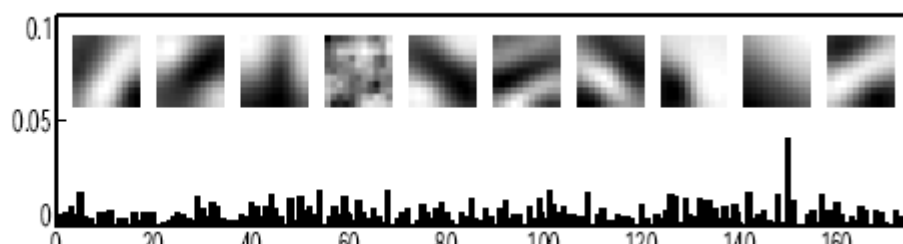
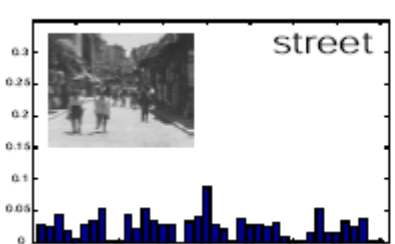
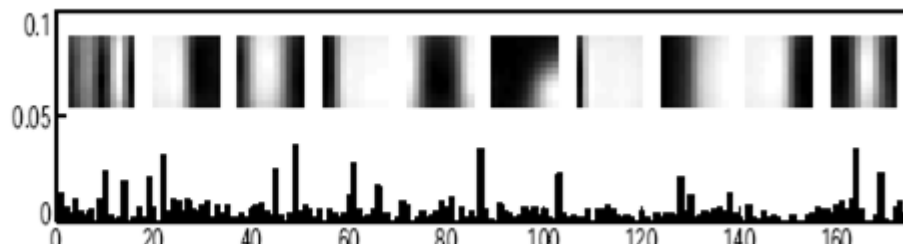
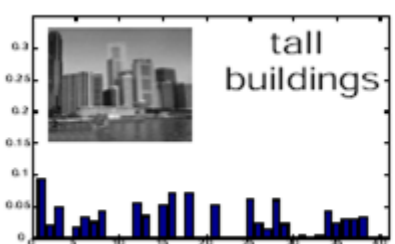
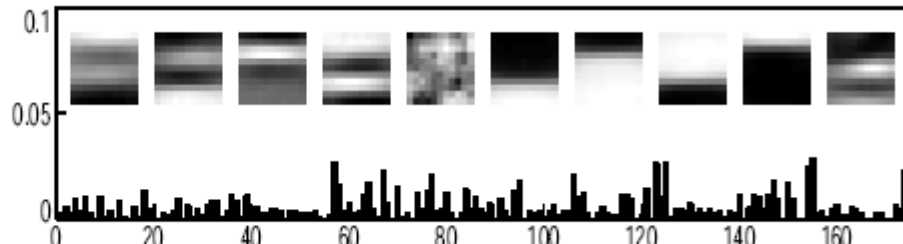
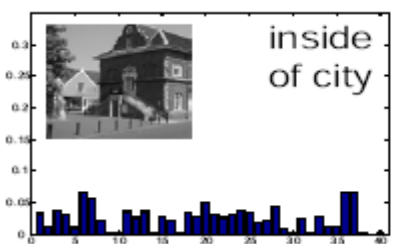
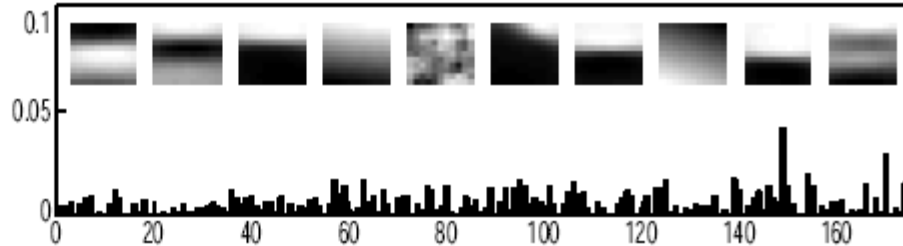
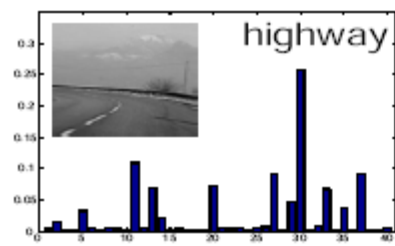


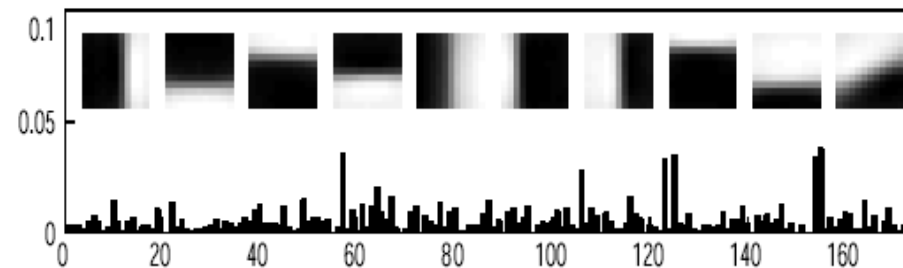
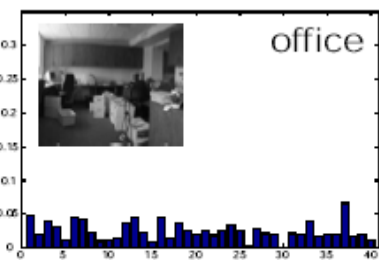
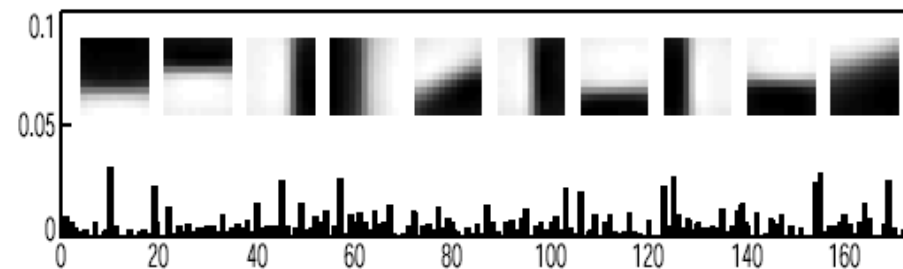
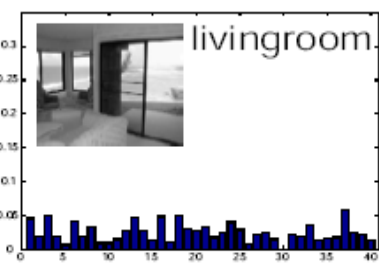
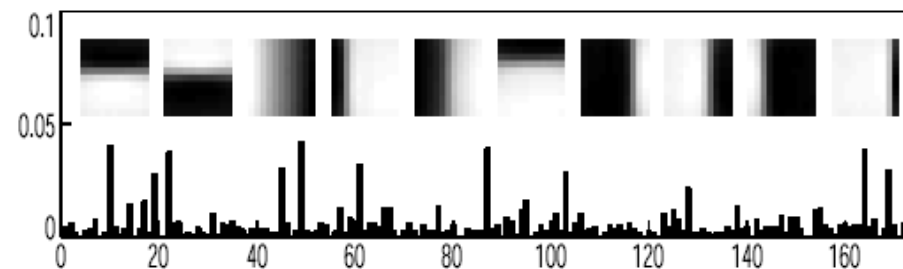
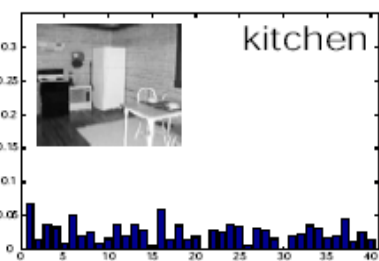
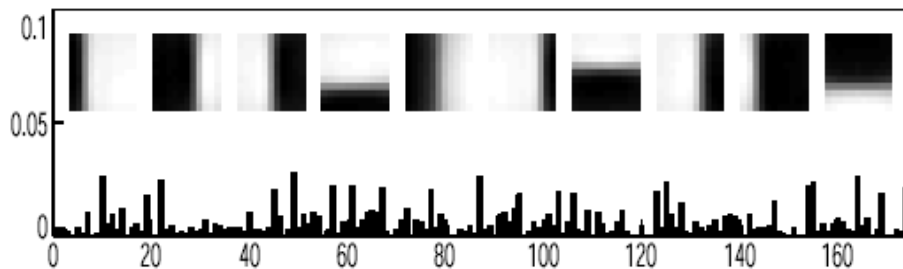
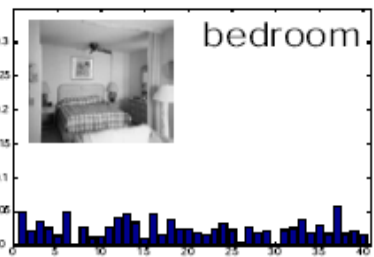
$$c = \arg \max_c p(x|c, \theta, \beta)$$

**category
decision**

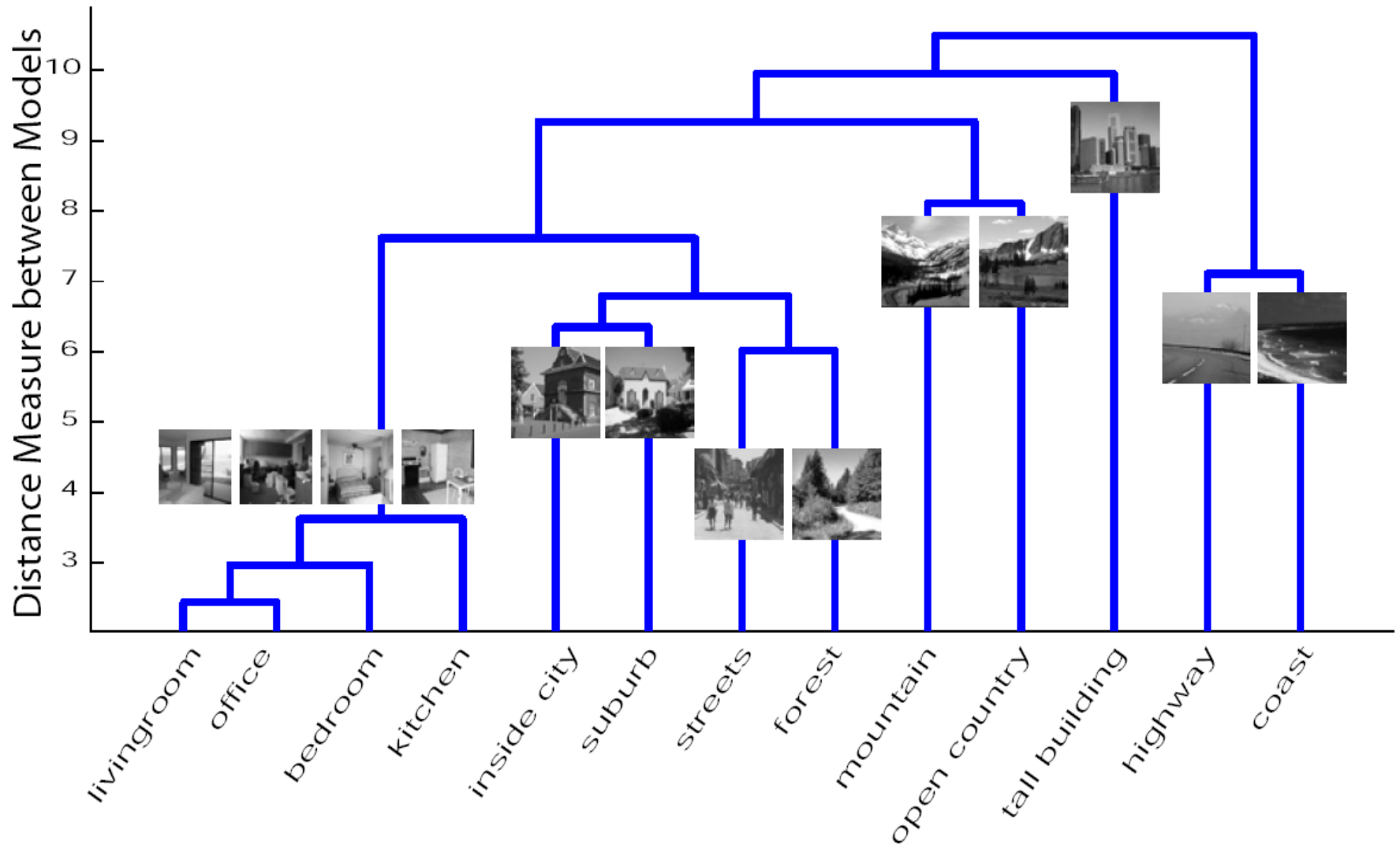
	highway	insidecity	tallbuildings	street	suburb	forest	coast	mountain	opencountry	bedroom	kitchen	livingroom	office
highway	74	2		2	2		14	4		2			
insidecity		58	10	6	8		4			2	6	4	2
tallbuildings		4	76	10				4		4		2	
street	2	4	6	78		2		2	2			4	
suburb					94					2			4
forest						88		12					
coast	2						78		20				
mountain	4		4		2	6	8	70	6				
opencountry	8				8	10	16	10	48				
bedroom	4	2	2		2	2	2	4		28	12	38	4
kitchen		8	2				2				60	14	14
livingroom		2	2	2			2	4		4	18	56	10
office					2		2			8	12	12	64







model distance based on topic distribution



Beaches



Highways



Forests



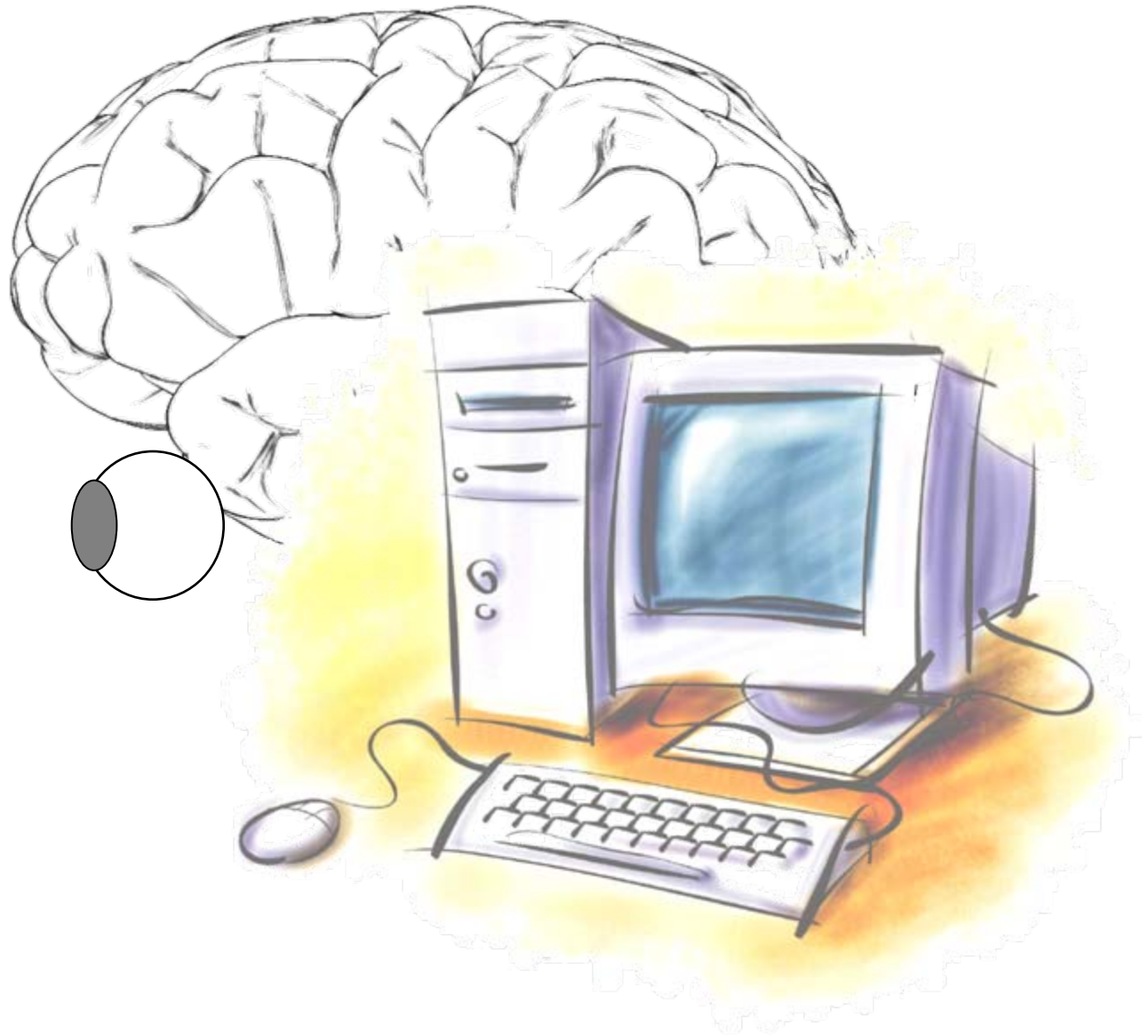
Buildings



Mountains



Industry



Change blindness



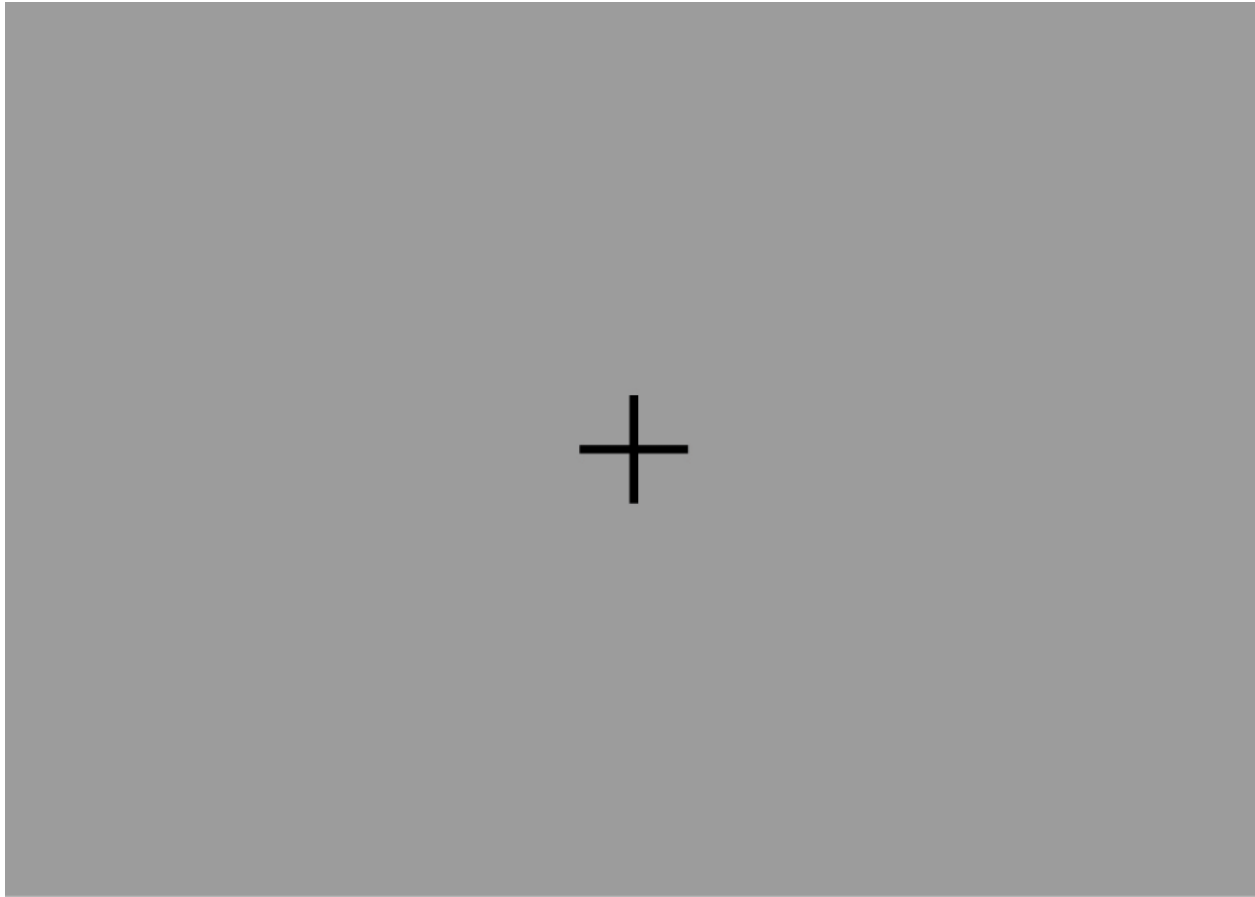
Rensink, O'regan, Simon, etc.

Change blindness



Rensink, O'regan, Simons, etc.

what **DO** we see in a glance?





Stage I: Collect Image Description

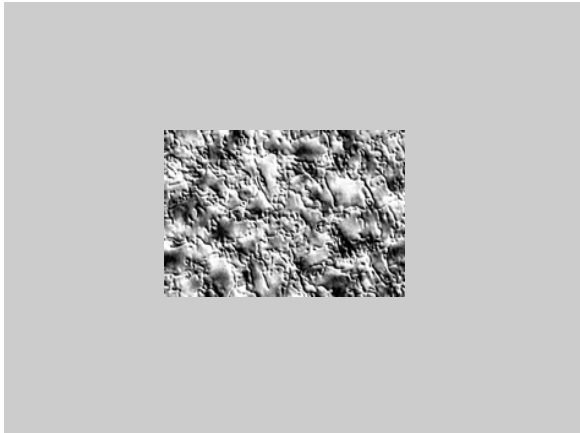
--- Illustration of 1 Trial

Subject types freely what he/she saw in the image



Please type your description here:

An outdoor scene, I think. reminded me a city... like walking in a park in new york or something. there seemed to be trees and a road and then this large skyscraper in the background.



time

Mask onset: $t = PT$



1 of 7 possible PT's (msec):
27, 40, 53, 67, 80, 120, 500

Image onset: $t = 0$ msec



PT = 500ms

This is indoors. It's must be a rich person's house. There are many paintings on the wall. The largest painting might have a fireplace beneath it. I think the largest painting was that of a man standing erect. The room is richly decorated and it looks like one of the rooms in Mr. Darcy's house in the A&E movie *Pride and Prejudice*. Or maybe it more closely resembles one of the rooms where the one of the rooms in Huntington's house (at the Huntington).

PT = 27ms

Couldn't see much; it was mostly dark w/ some square things, maybe furniture. (Subject: AM)

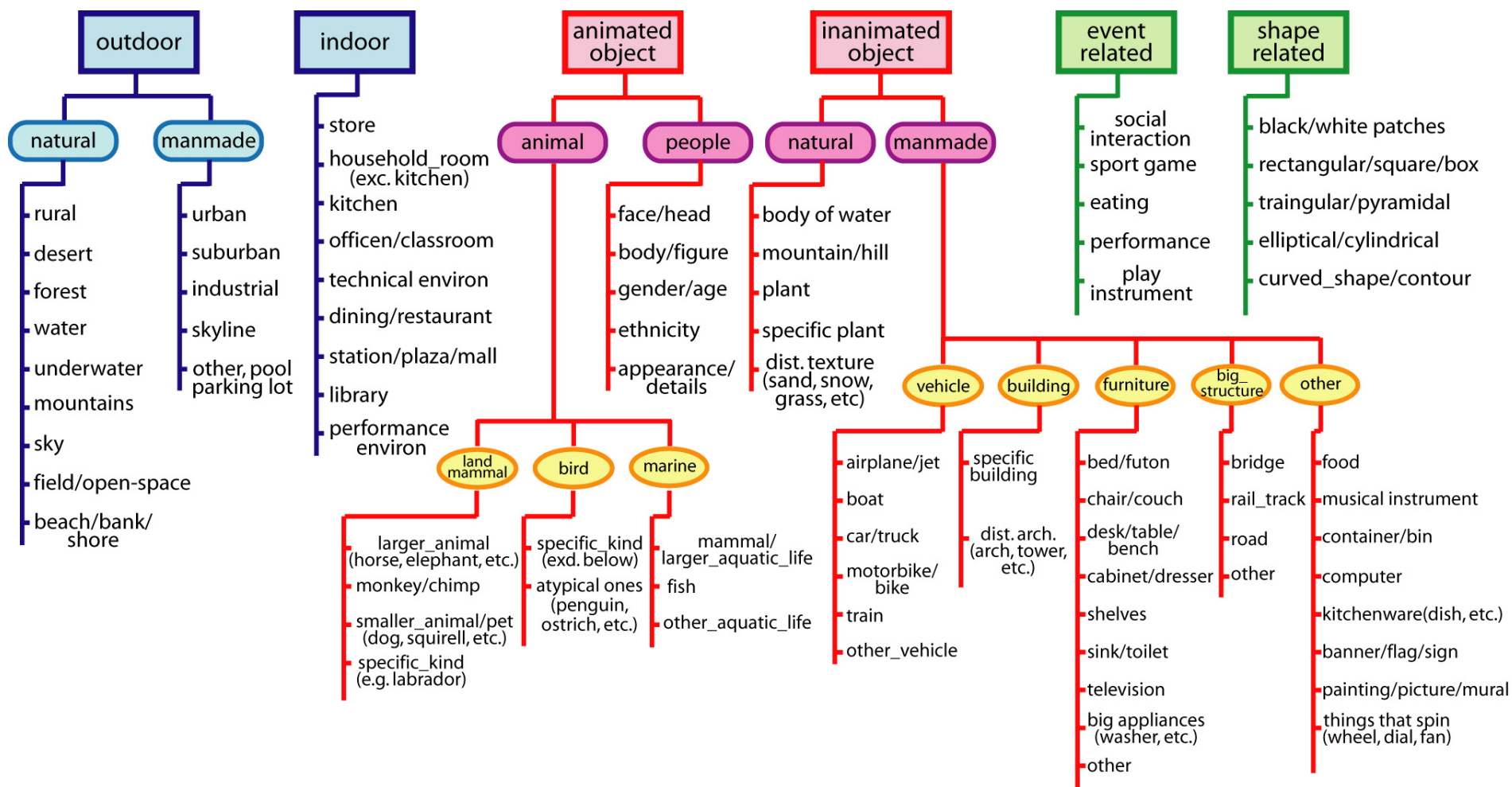
PT = 40ms

This looked like an indoor shot. Saw what looked like a large framed object (a painting?) on a white background (i.e., the wall). (Subject: RW)

PT = 67ms

I saw the interior of a room in a house. There was a picture to the right, that was black, and possibly a table in the center. It seemed like a formal dining room. (Subject: JB)

Response attributes??



Response No. 18 for Image No. 4



I could make out some kind of circular shapes near the bottom of the picture. These reminded me of those round life preservers that are on ships. There was also a man standing on top of some wooden structure.

CATEGORY: SENSORY/SHAPES

Please select one of "correct" or "incorrect" for each checked description. Click "Next>>" to continue

- | | | | |
|---------------------------------------|---|--|---------------------------------|
| black/white_patches | <input type="checkbox"/> described | <input checked="" type="radio"/> correct | <input type="radio"/> incorrect |
| rectangular/square/box | <input type="checkbox"/> described | <input checked="" type="radio"/> correct | <input type="radio"/> incorrect |
| triangular/pyramidal | <input type="checkbox"/> described | <input checked="" type="radio"/> correct | <input type="radio"/> incorrect |
| elliptical/cylindrical(eg.round,blob) | <input checked="" type="checkbox"/> described | <input checked="" type="radio"/> correct | <input type="radio"/> incorrect |
| curved_shape/contour(eg.arc,'S') | <input type="checkbox"/> described | <input checked="" type="radio"/> correct | <input type="radio"/> incorrect |

What's in a glance?

Average fixation time (one glance) = 120-200ms



What's in a glance?



PT = 107ms

This is outdoors. A black, furry dog is running/walking towards the right of the picture. His tail is in the air and his mouth is open. Either he had a ball in his mouth or he was chasing after a ball. (Subject EC)

PT = 500ms

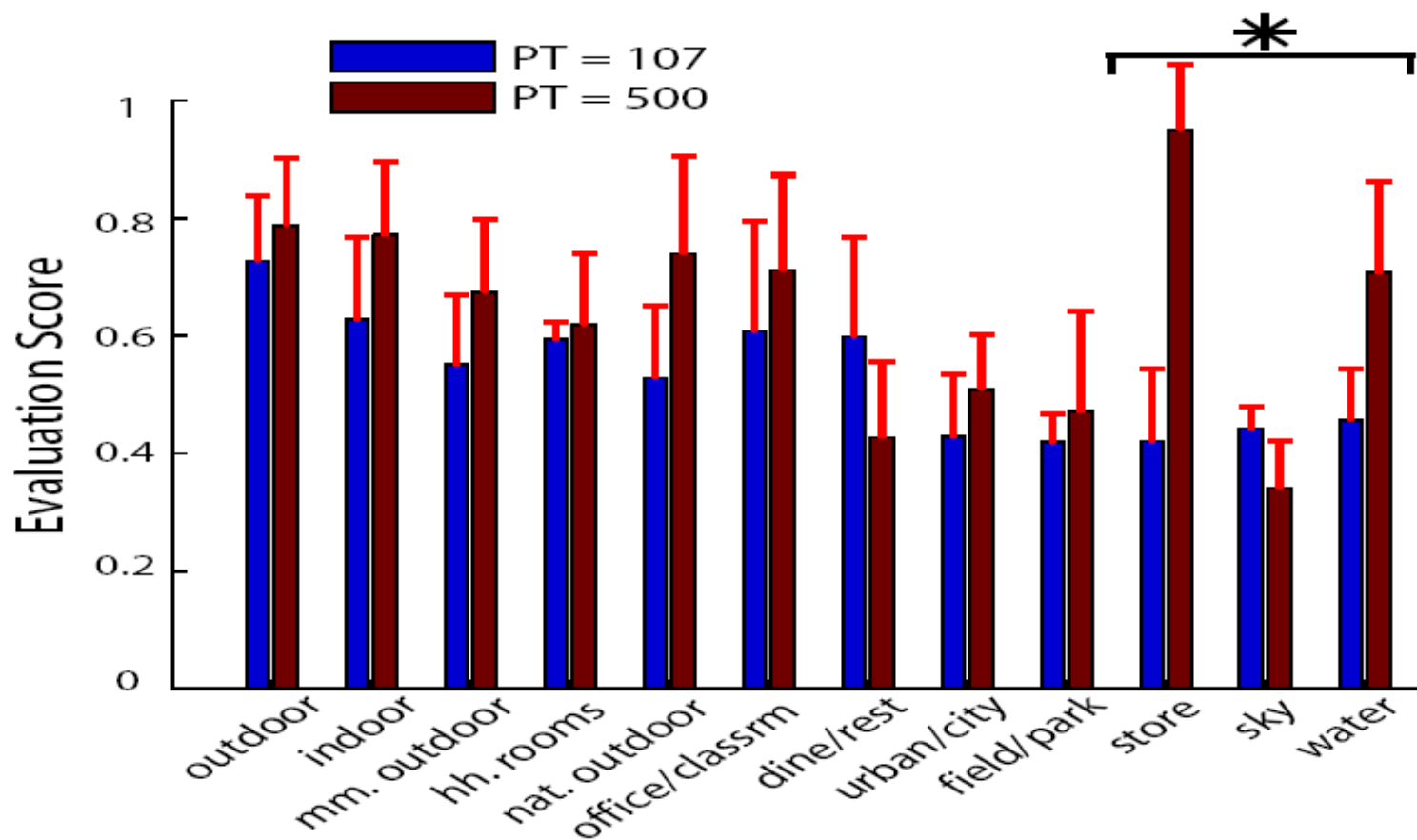
I saw a black dog carrying a gray frisbee in the center of the photograph. The dog was walking near the ocean, with waves lapping up on the shore. It seemed to be a gray day out. (Subject JB)



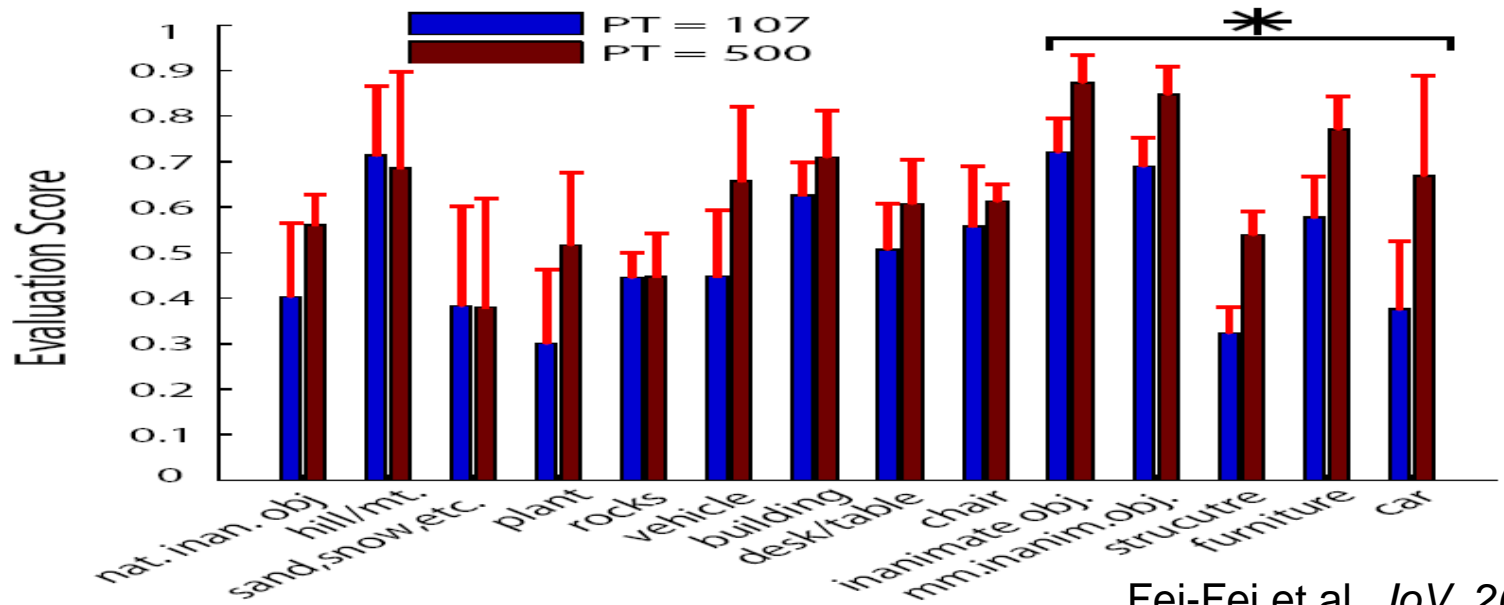
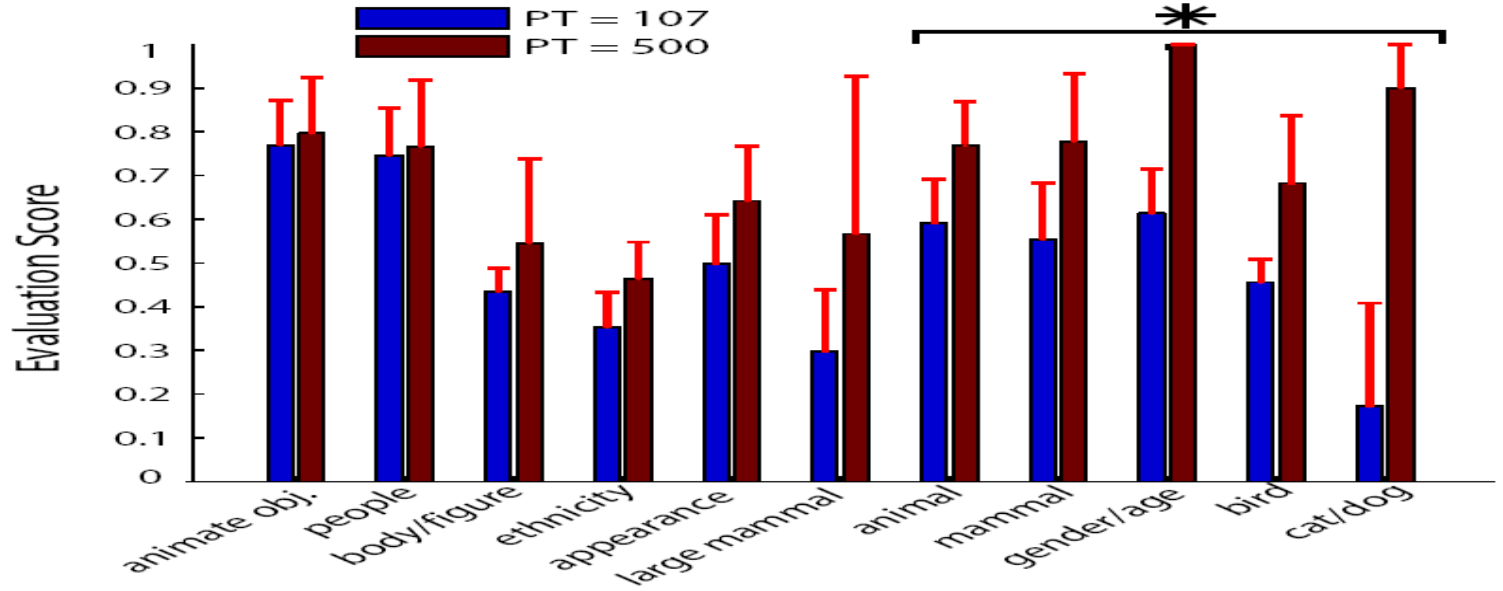
inside a house, like a living room, with chairs and sofas and tables, no ppl. (Subject HS)

A room full of musical instruments. A piano in the foreground, a harp behind that, a guitar hanging on the wall (to the right). It looked like there was also a window behind the harp, and perhaps a bookcase on the left. (Subject RW)

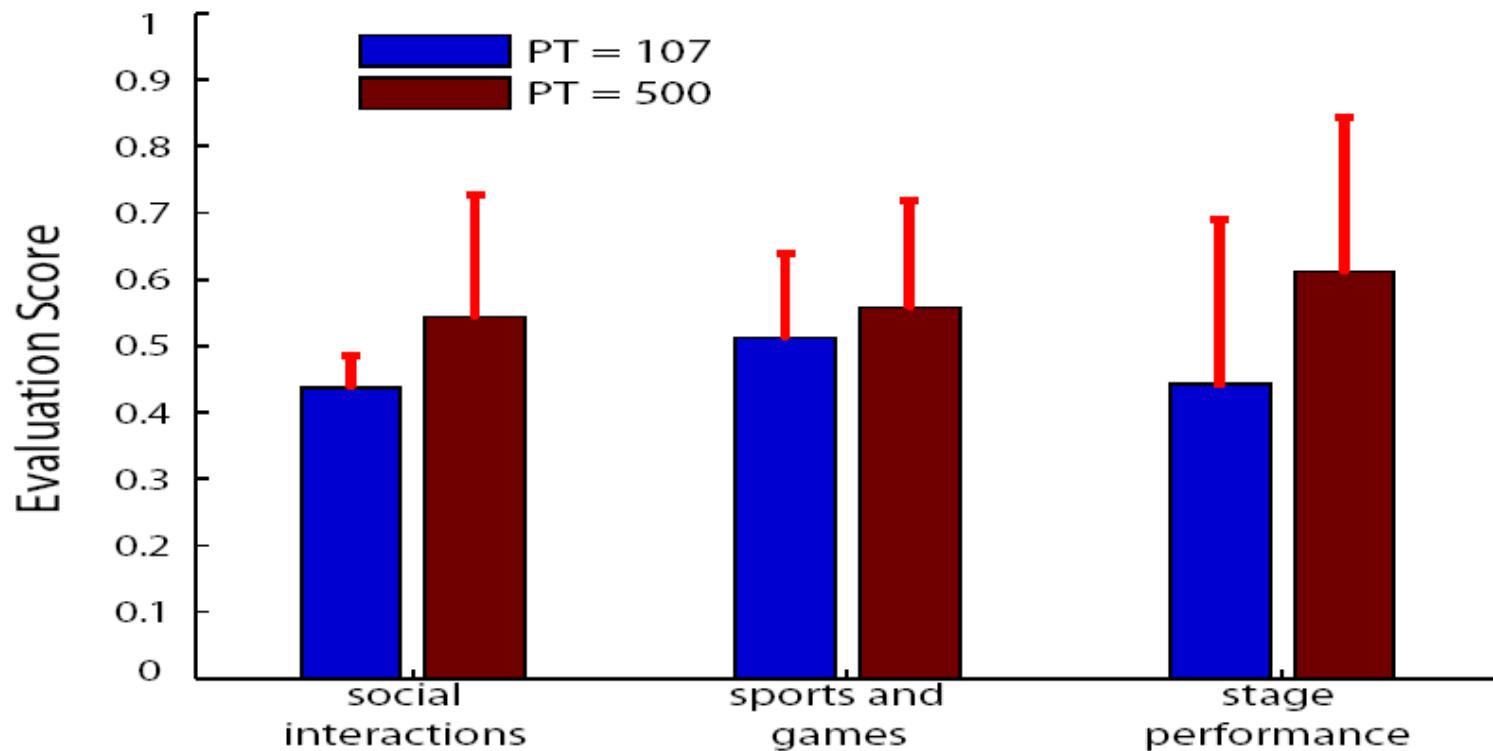
Scene level



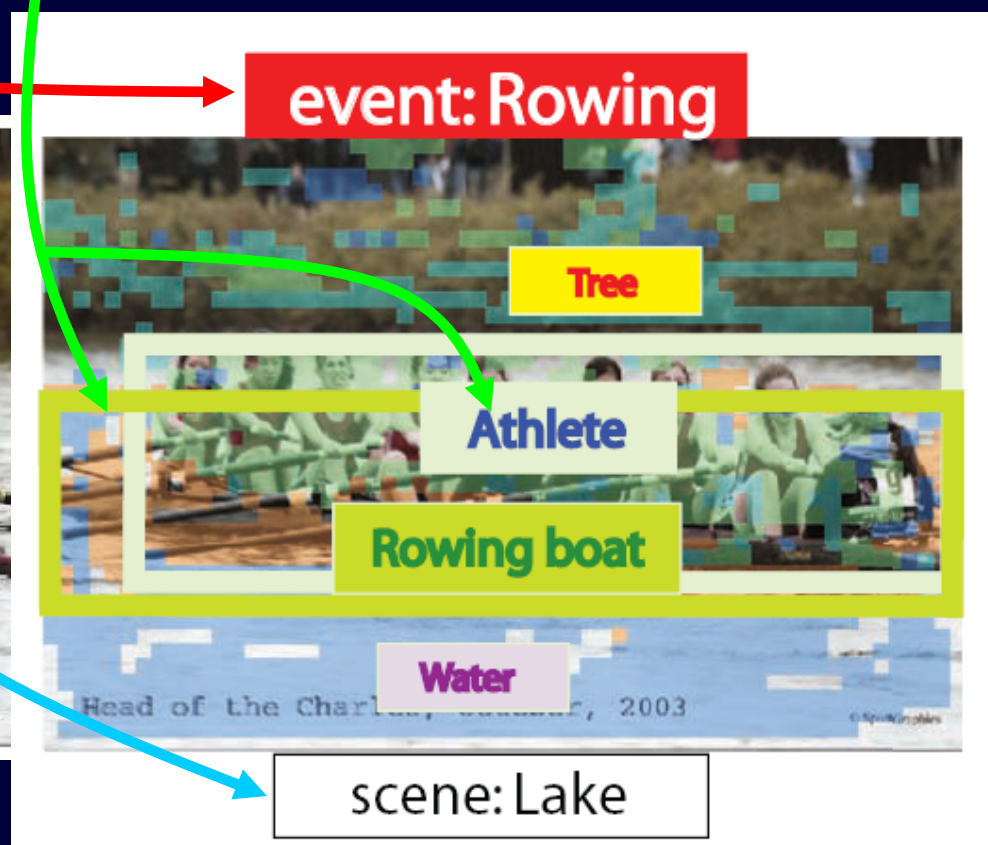
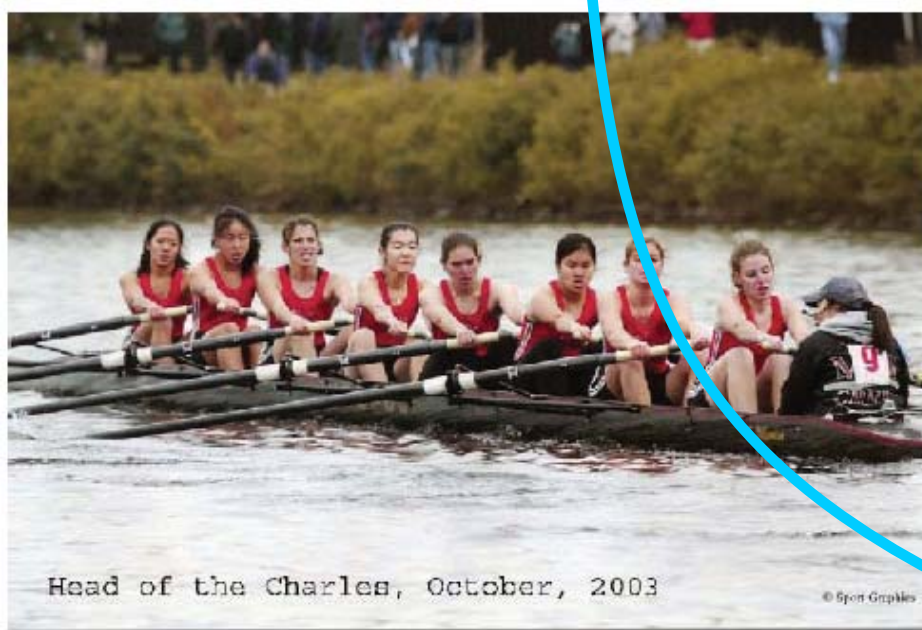
Object level

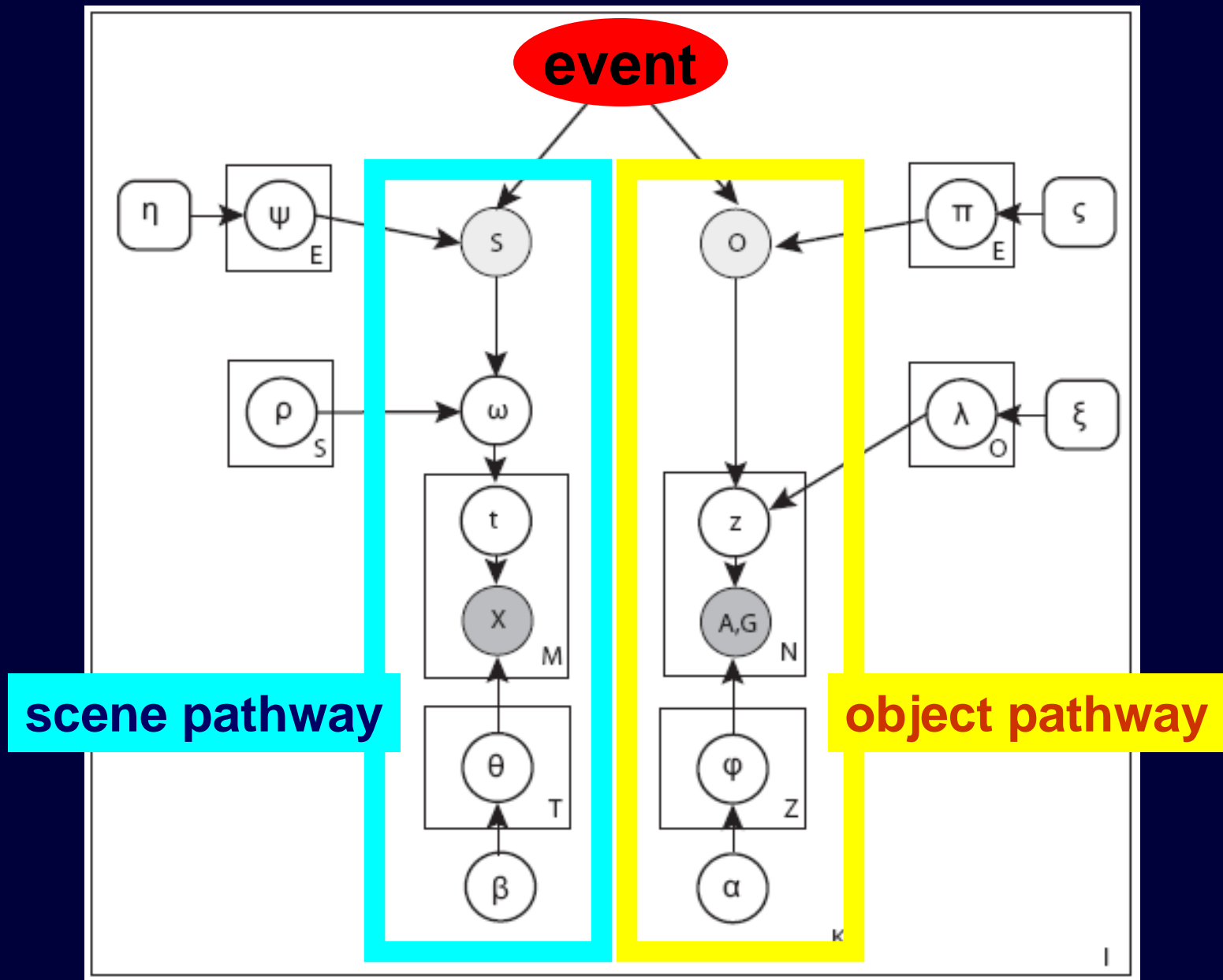


(Social) Events



What, where and who? Classifying events by scene and object recognition

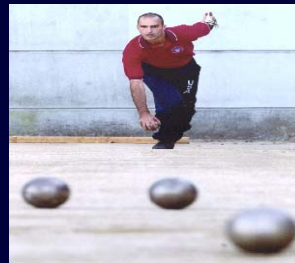




rowing



bocce



badminton



snow boarding



polo



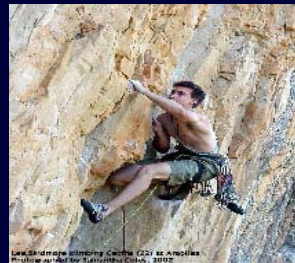
croquet



sailing



rock climbing



event: Badminton



scene: Badminton court

event: Bocce



scene: Bocce court

event: Croquet

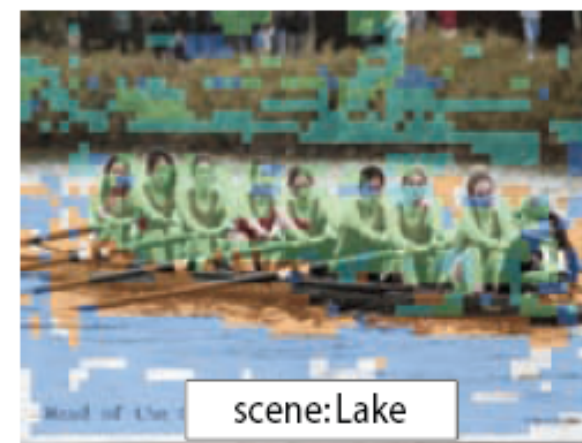
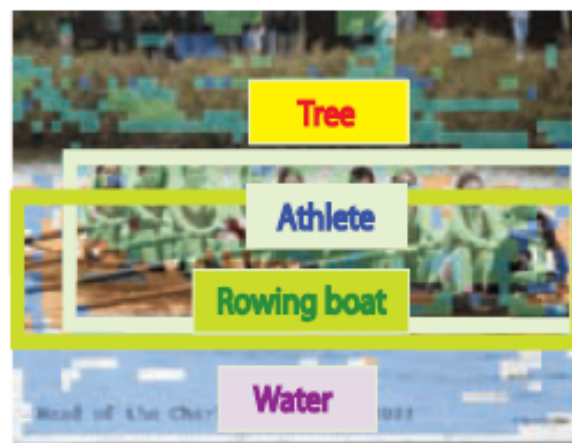


scene: Croquet court

event: Rockclimbing



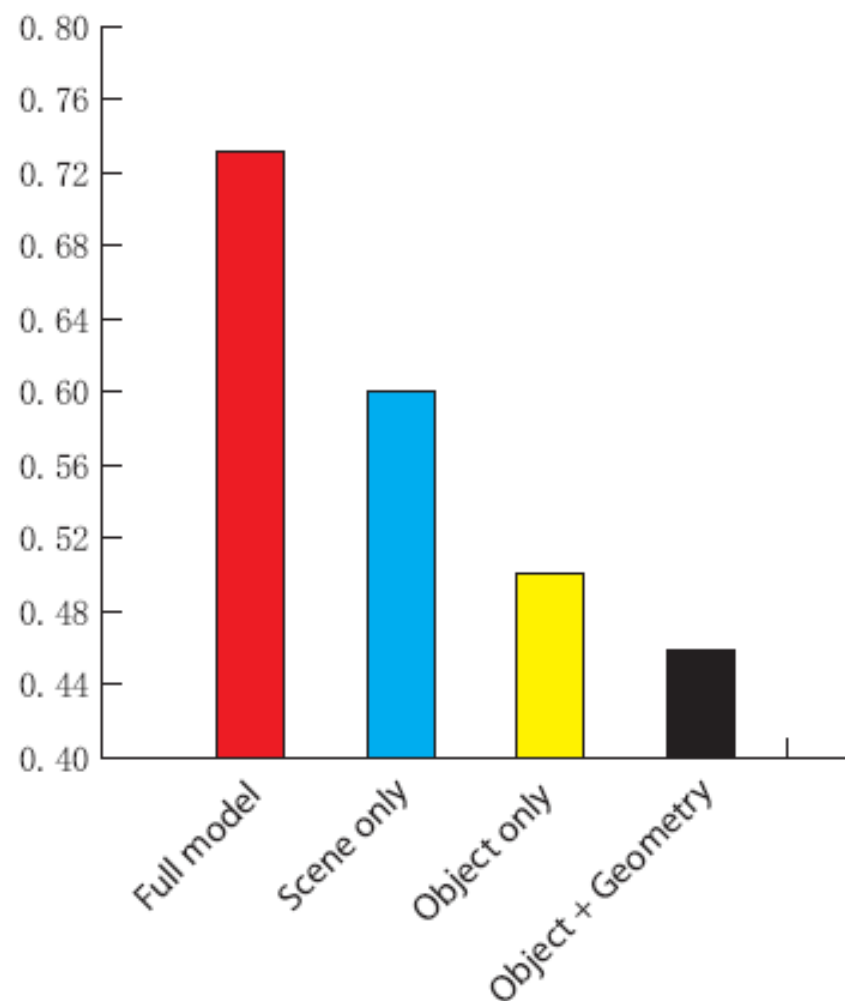
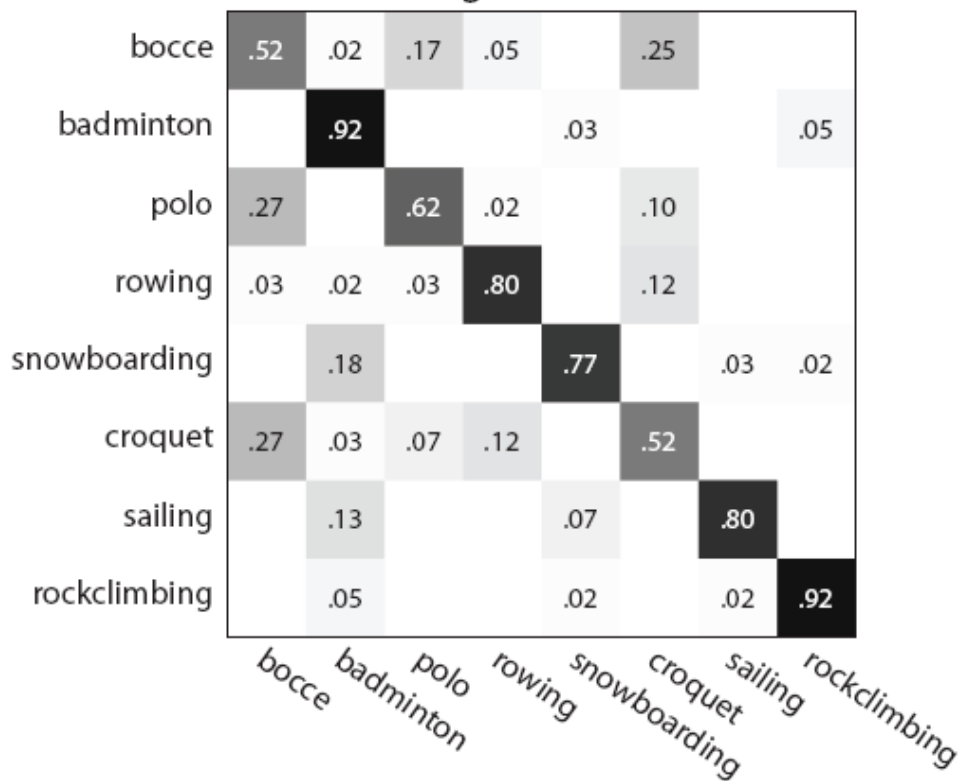
event: Rowing



event: Sailing

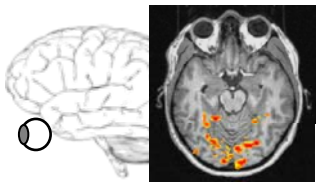


Average Perf. = 73.4%

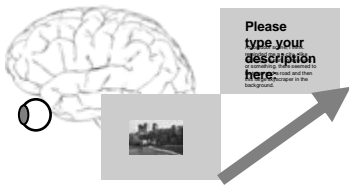




#1: natural scene categorization entails little attention



#2: decoding the neural representation of natural scene categories



#3: what can we perceive within a glance of a scene?



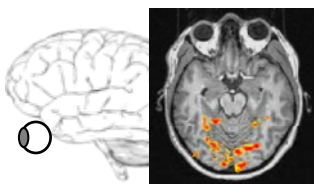
#4: Bayesian graphical models for natural scene categorization and event recognition



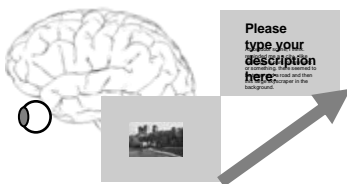
Thank you!



#1: natural scene categorization entails little attention (Rufin VanRullen, Pietro Perona, Christof Koch)



#2: decoding the neural representation of natural scene categories (Eamon Caddigan, Dirk Walther, Diane Beck)



#3: what can we perceive within a glance of a scene? (Asha Iyer, Pietro Perona, Christof Koch)



#4: Bayesian graphical models for natural scene categorization and event recognition (Pietro Perona, Li-Jia Li)

