

Using the Forest to see the Trees: A computational model relating features, objects and scenes

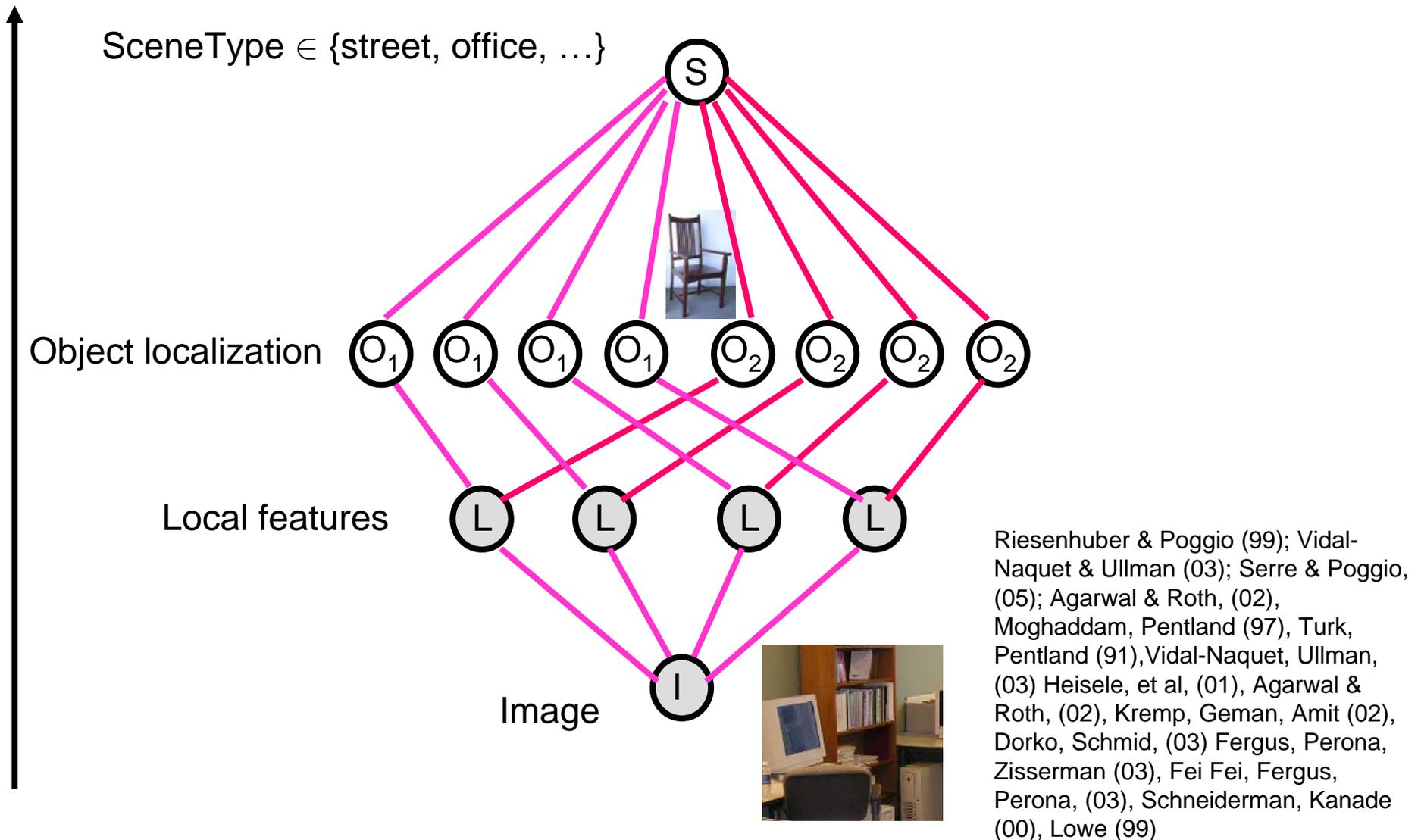
Antonio Torralba
CSAIL-MIT

Joint work with

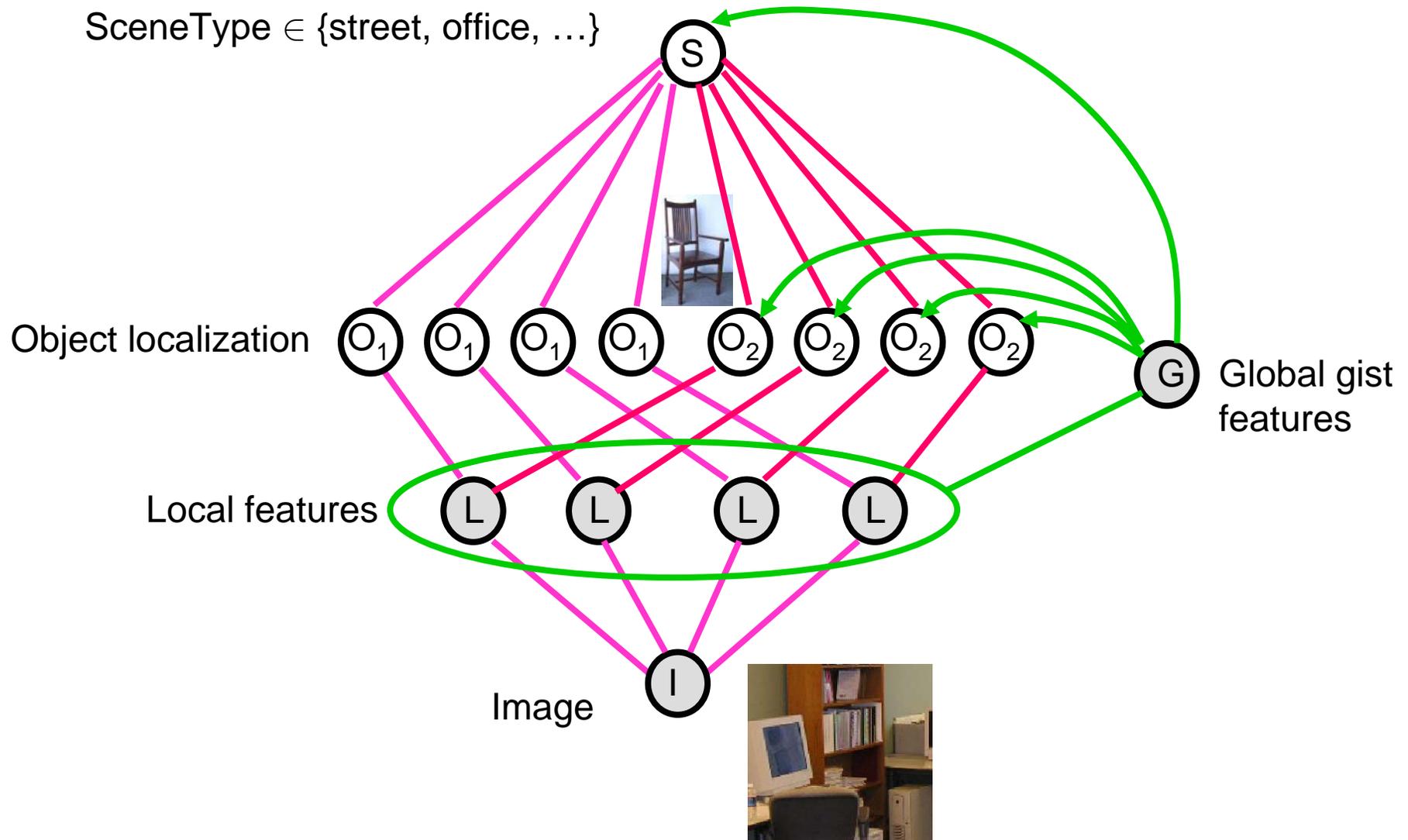
Aude Oliva, Kevin Murphy, William Freeman

Monica Castelhana, John Henderson

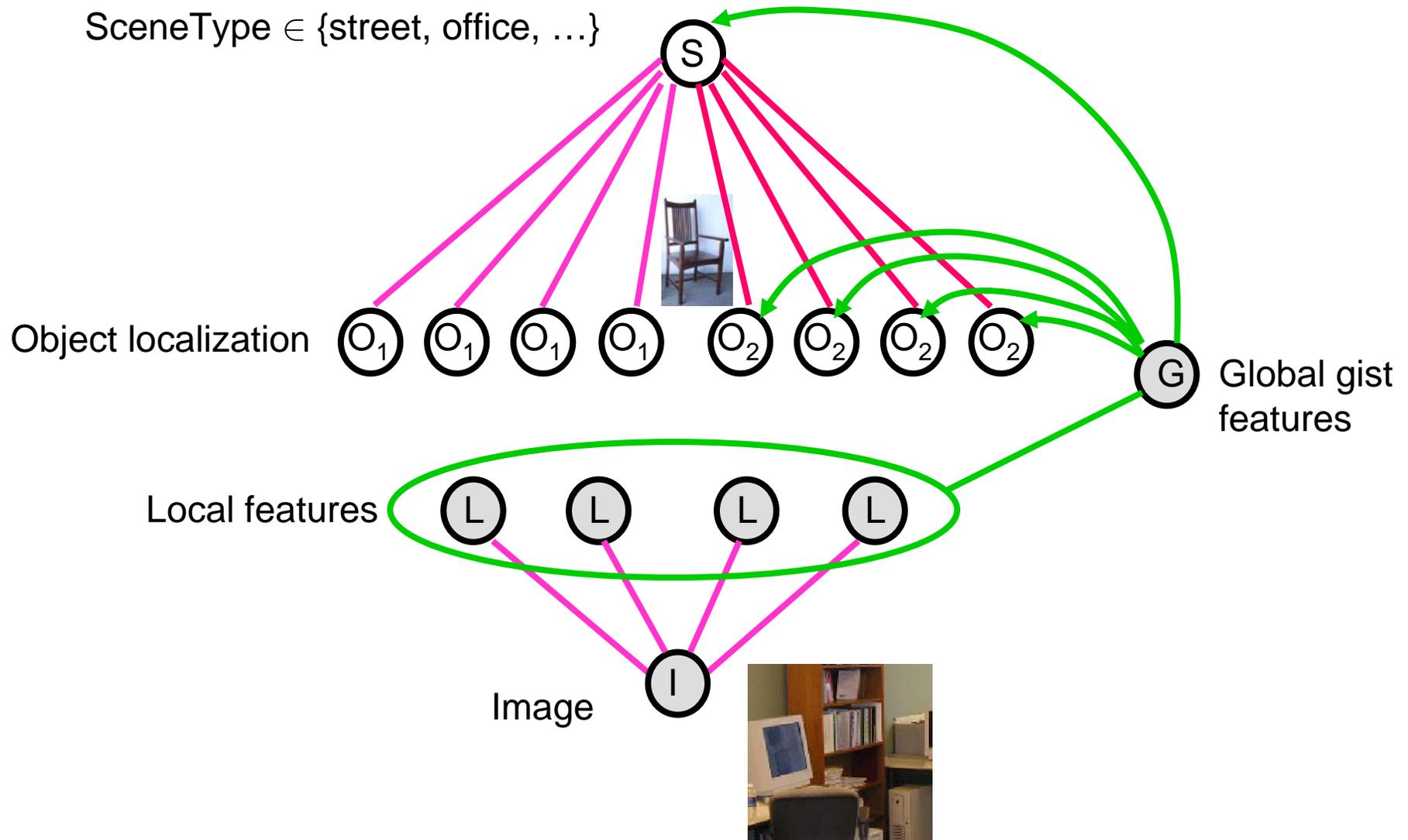
From objects to scenes



From scenes to objects



From scenes to objects



The context challenge

What do you think are the hidden objects?



The context challenge

What do you think are the hidden objects?

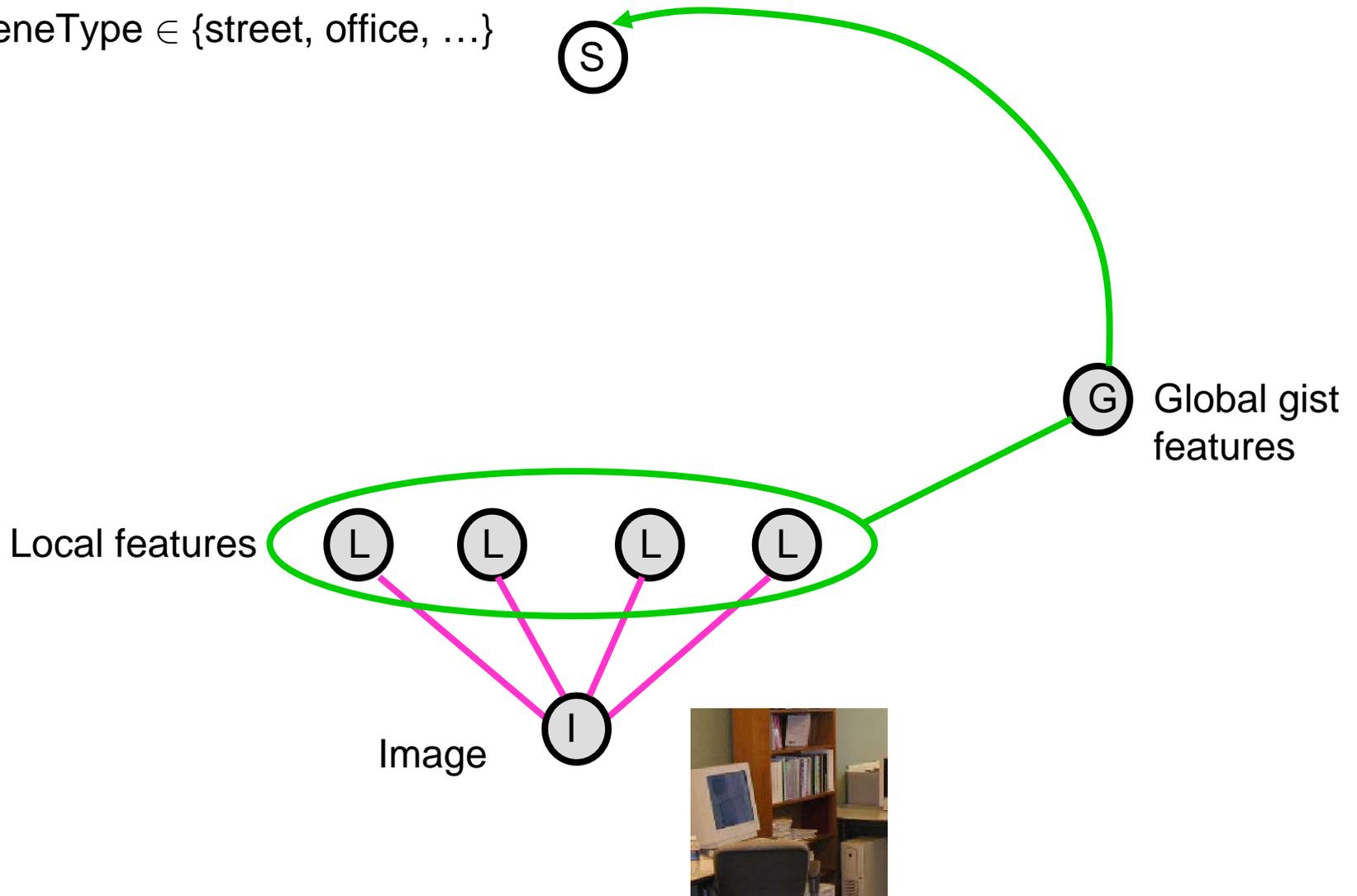


Chance ~ 1/30000

Answering this question does not require knowing how the objects look like. It is all about context.

From scenes to objects

SceneType \in {street, office, ...}



Scene categorization

Office



Corridor



Street



Place identification

Office 610



Office 615



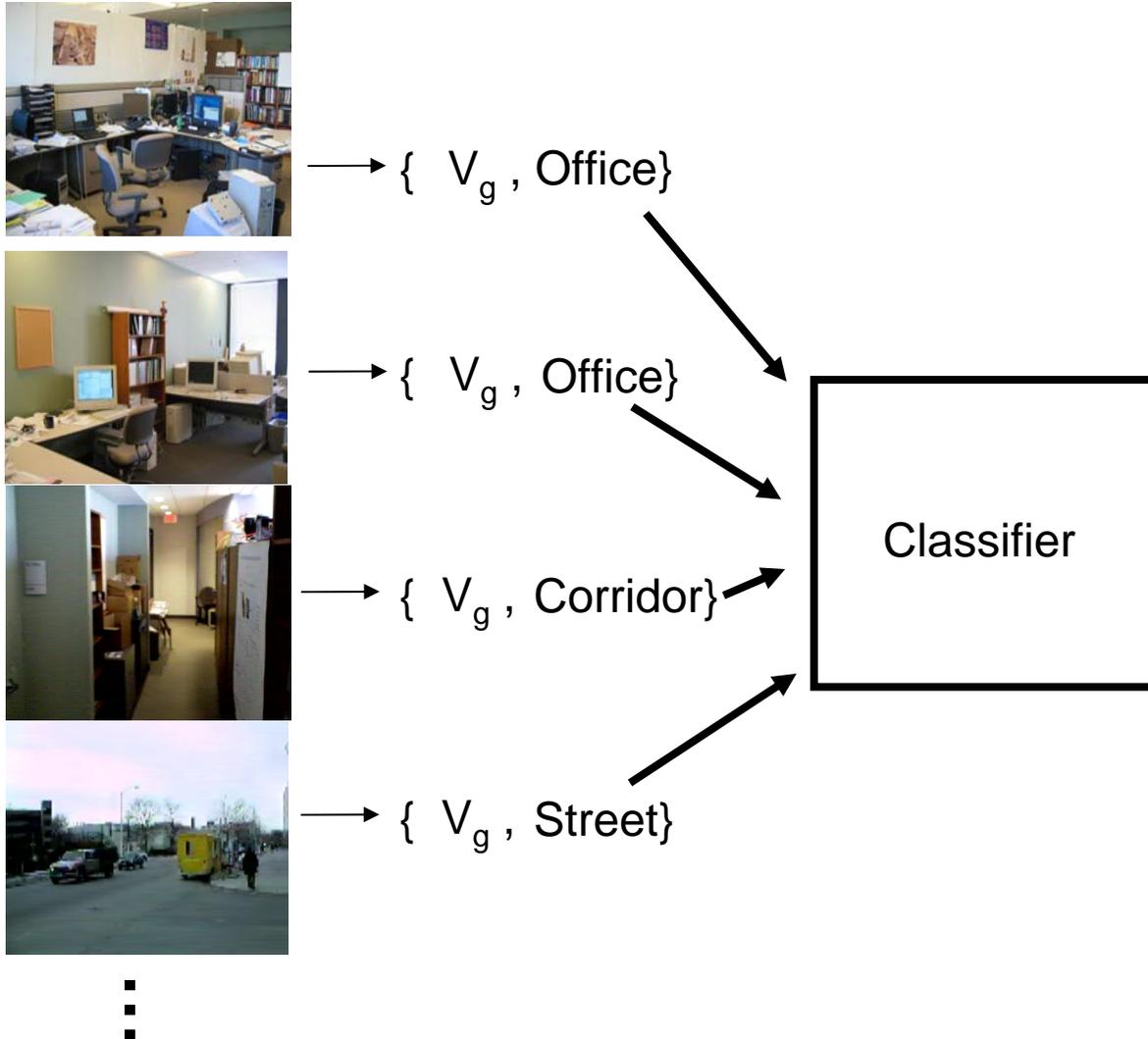
Draper street



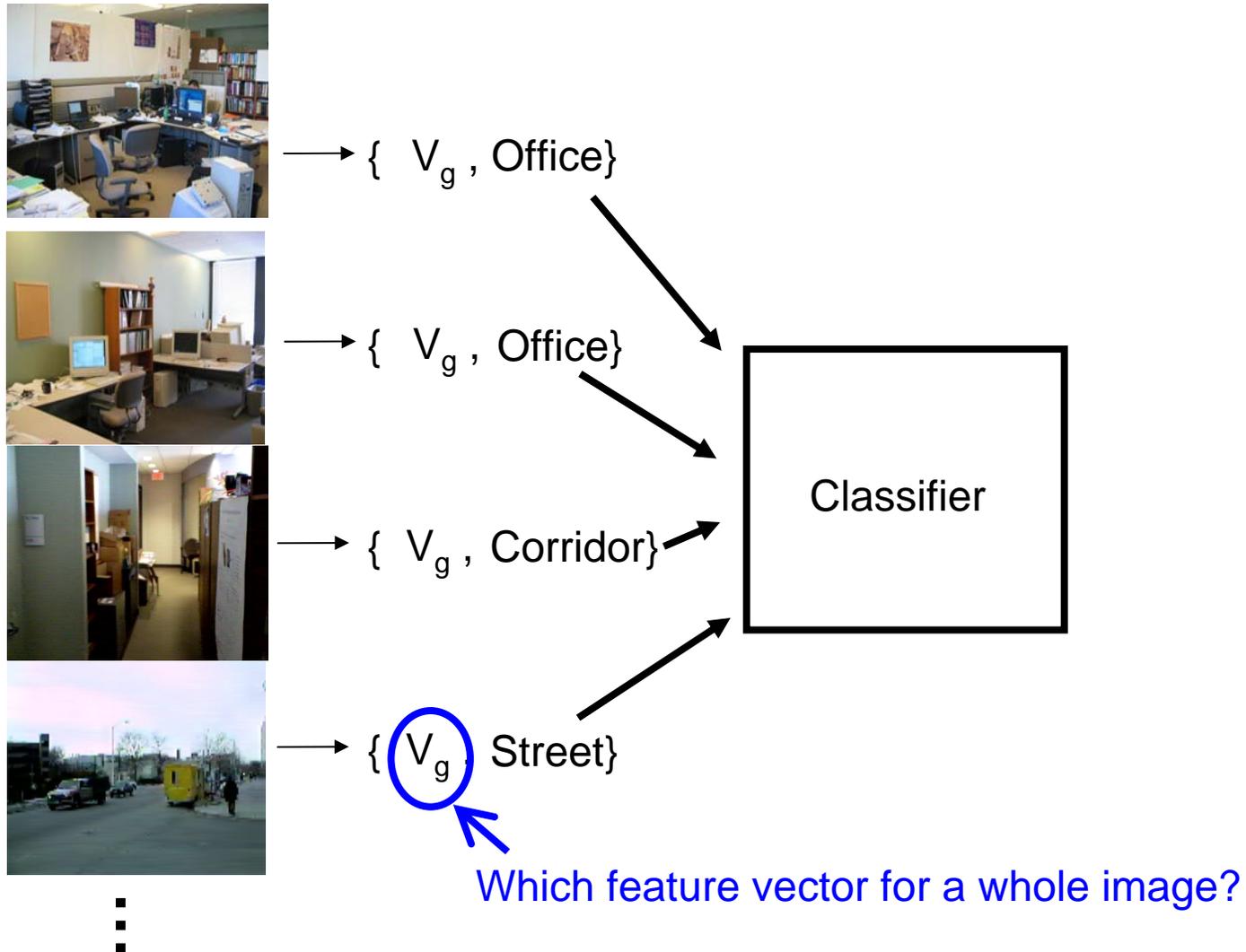
59 other places...

Scenes are categories, places are instances

Supervised learning

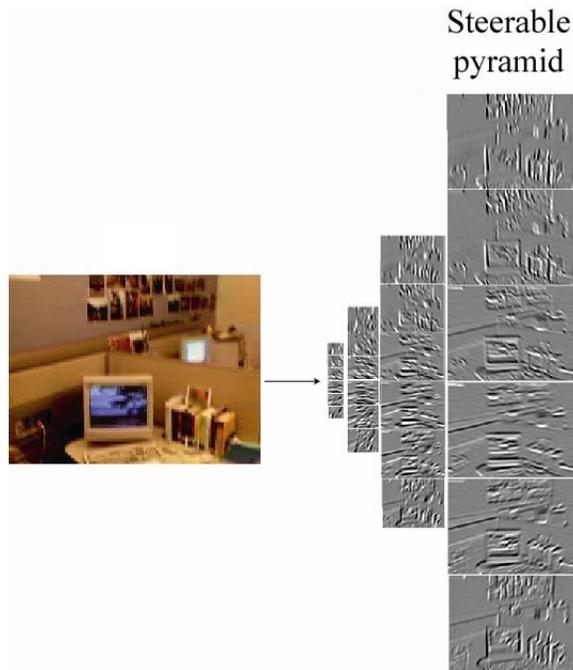


Supervised learning



Global features (gist)

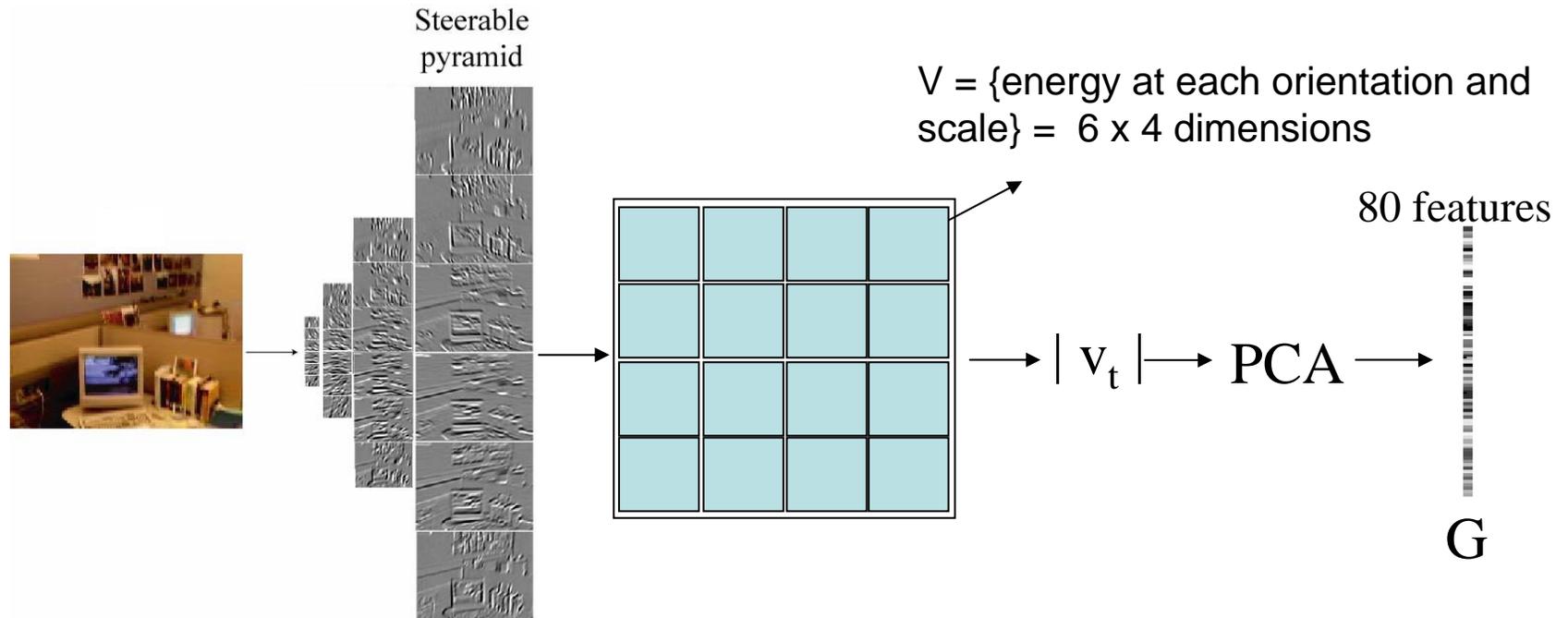
First, we propose a set of features that do not encode specific object information



Oliva & Torralba, IJCV'01; Torralba, Murphy, Freeman, Mark, CVPR 03.

Global features (gist)

First, we propose a set of features that do not encode specific object information

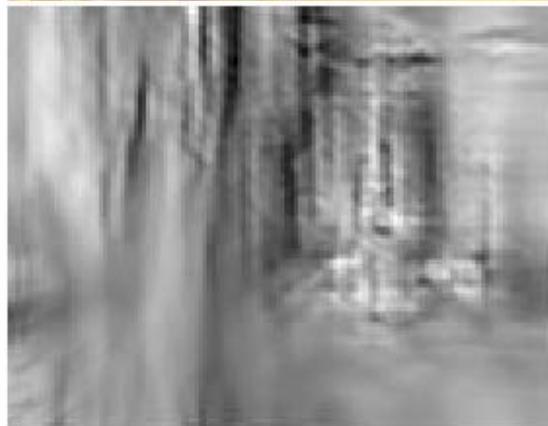


Example visual gists

I



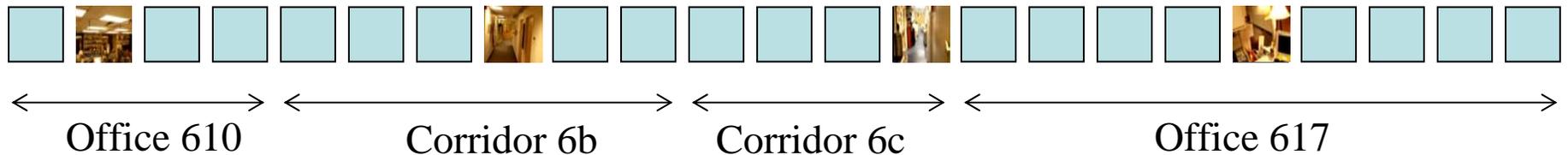
I'



Global features (I) ~ global features (I')

Learning to recognize places

We use annotated sequences for training



- Hidden states = location (63 values)
- Observations = v_t^G (80 dimensions)
- Transition matrix encodes topology of environment
- Observation model is a mixture of Gaussians centered on prototypes (100 views per place)

Wearable test-bed v1

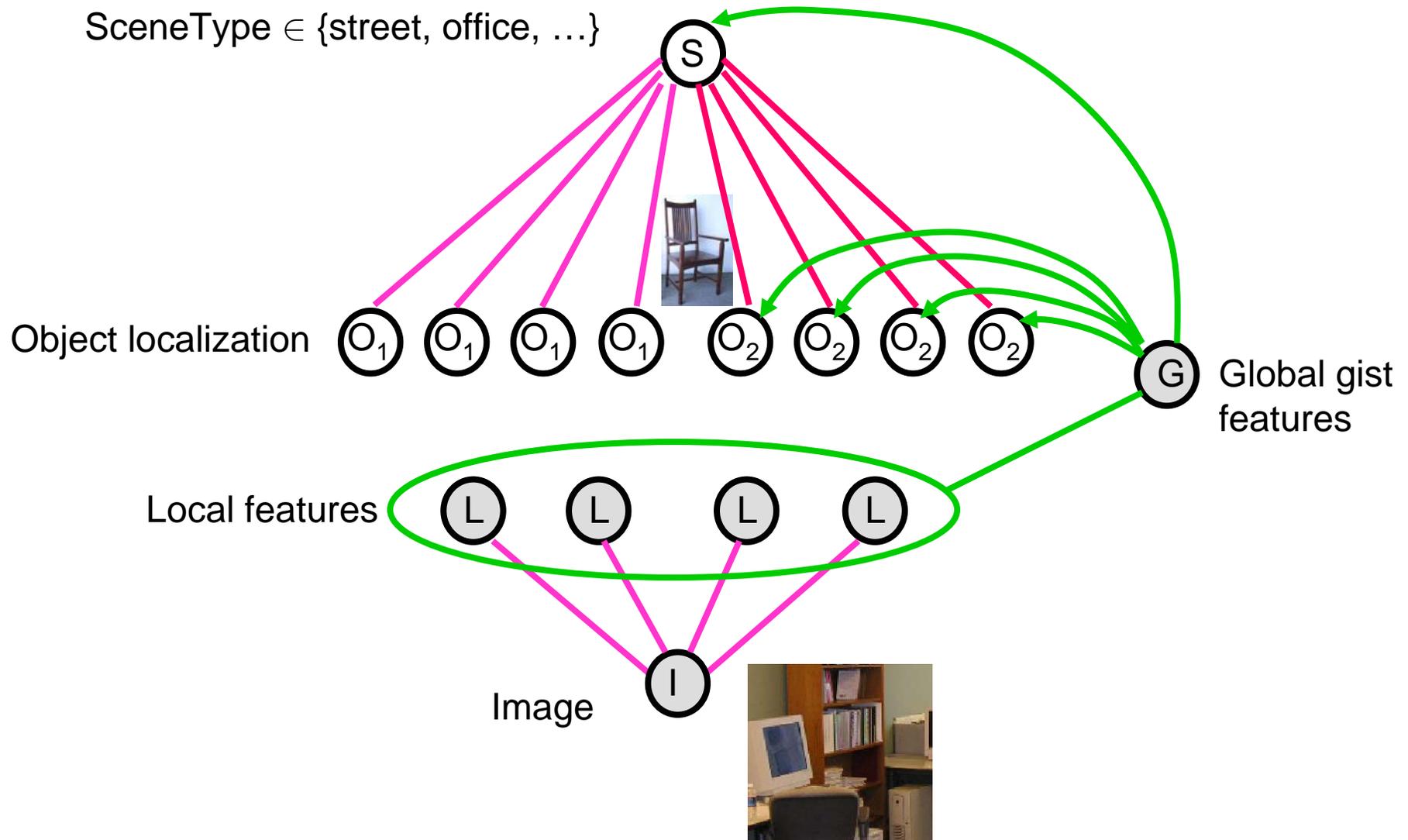
Wearable test-bed v2

Place/scene recognition demo

t=930, truth = 400-fl6-visionArea1



From scenes to objects



Global scene features predicts object location

New image

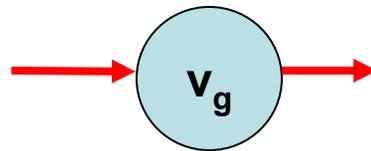
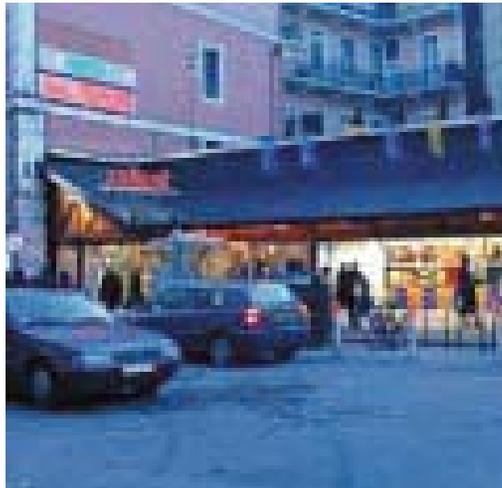


Image regions likely to contain the target

Global scene features predicts object location

Training set (cars)



→ $\{V_g^1, X^1\}$



→ $\{V_g^2, X^2\}$



→ $\{V_g^3, X^3\}$

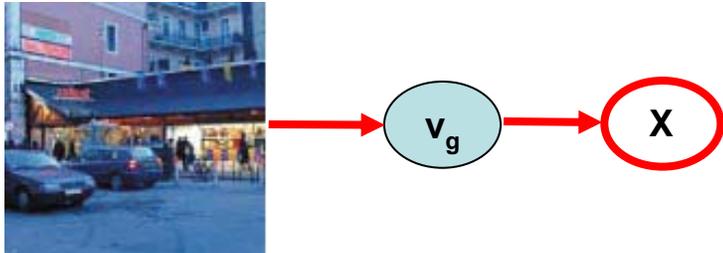


→ $\{V_g^4, X^4\}$

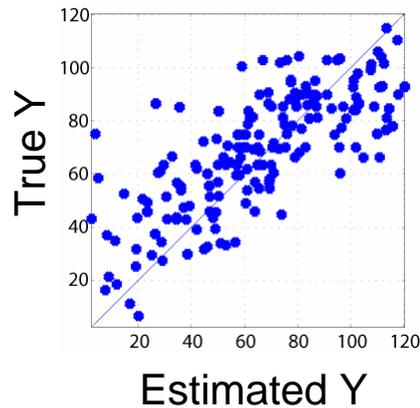
⋮

The goal of the training is to learn the association between the location of the target and the global scene features

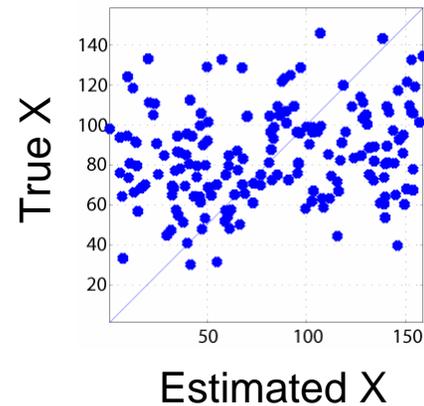
Global scene features predicts object location



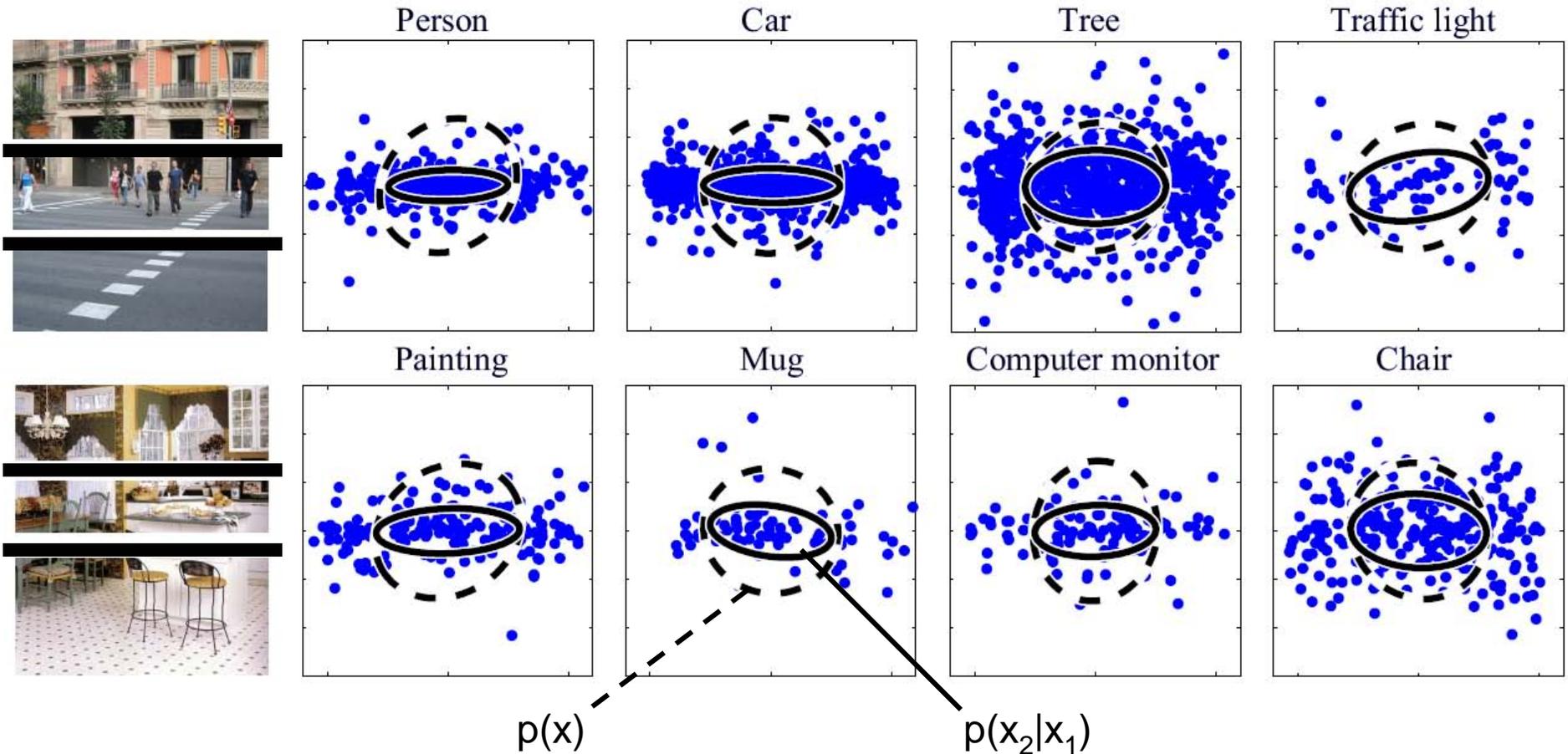
Results for predicting the vertical location of people



Results for predicting the horizontal location of people

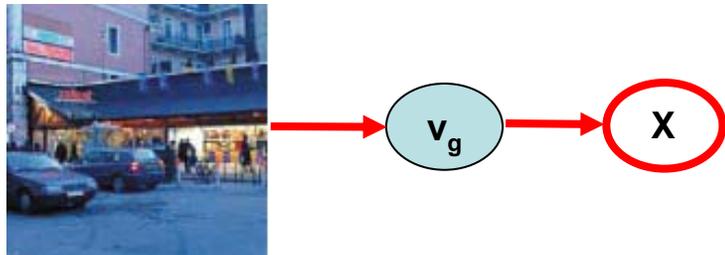


The layered structure of scenes



In a display with multiple targets present, the location of one target constraints the 'y' coordinate of the remaining targets, but not the 'x' coordinate.

Global scene features predicts object location



Stronger contextual constraints can be obtained using other objects.





SALE

\$15.00

WHITLOCK'S BAKERY

Whitlock's
BAKERY

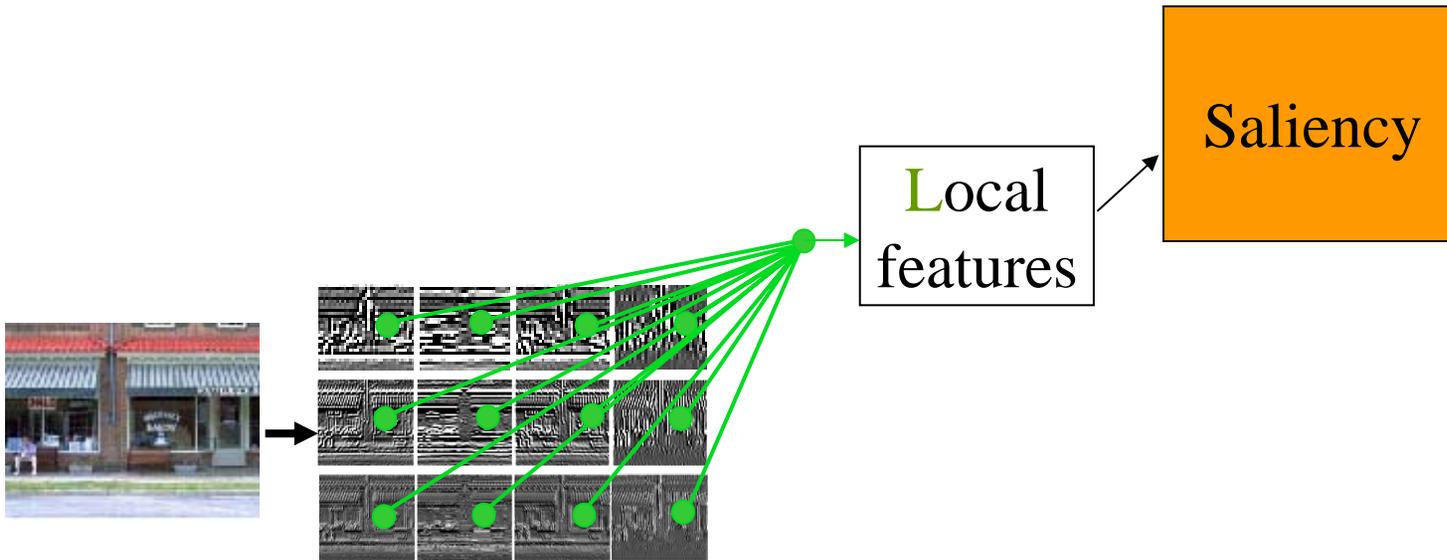
Whitlock's
BAKERY





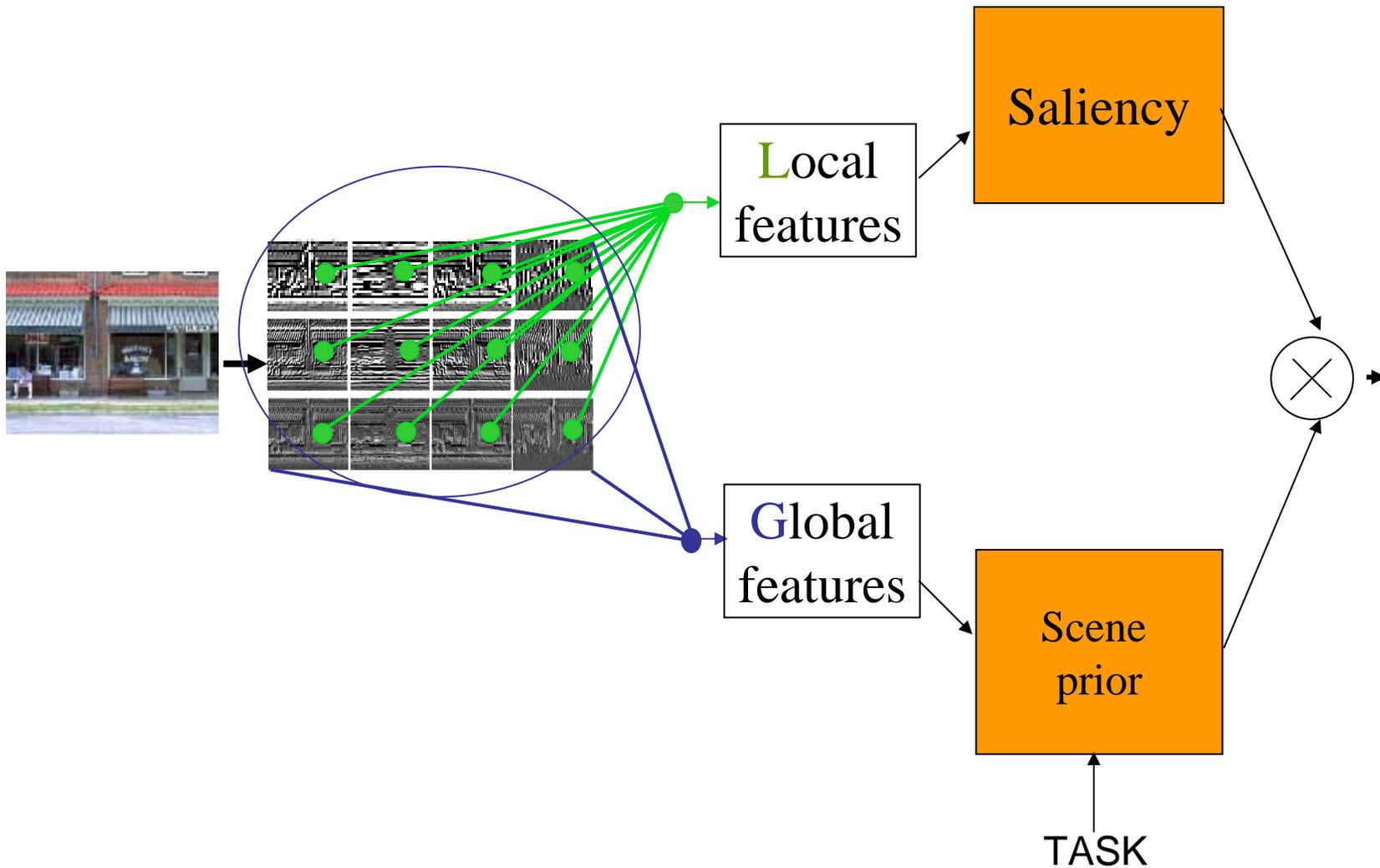
1

Attentional guidance

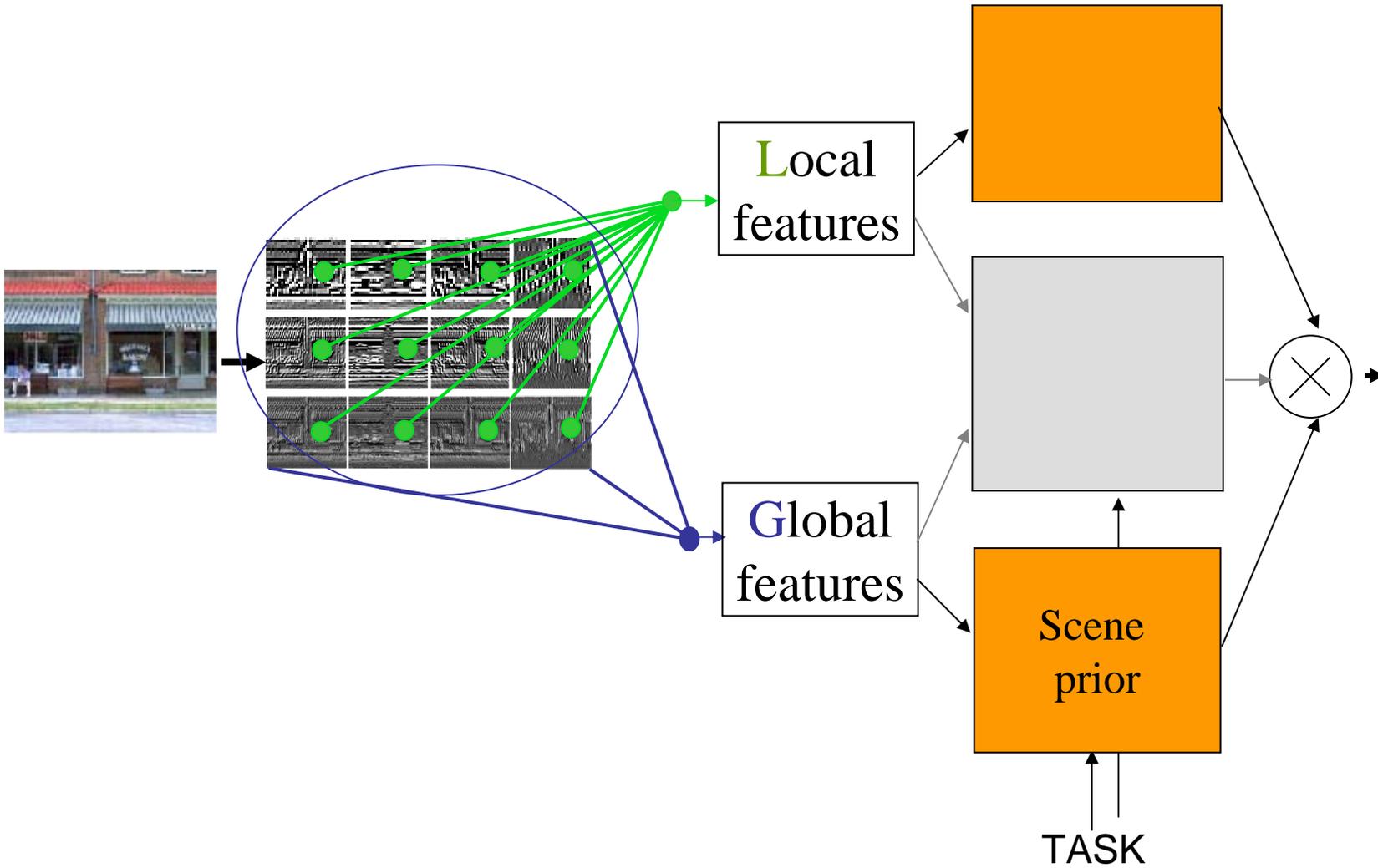


Saliency models: Koch & Ullman, 85; Wolfe 94; Itti, Koch, Niebur, 98; Rosenholtz, 99

Attentional guidance



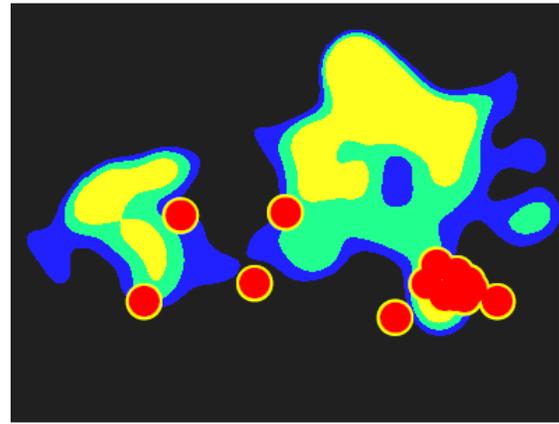
Attentional guidance



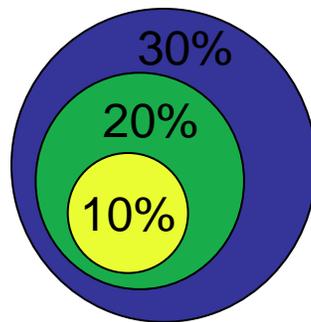
Comparison regions of interest



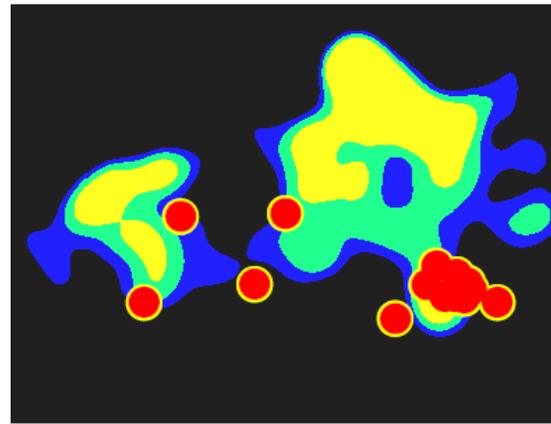
Comparison regions of interest



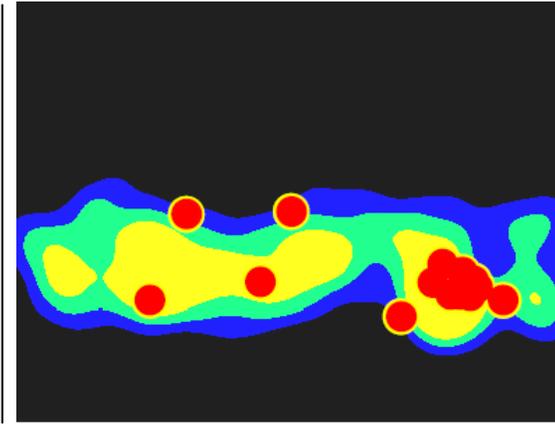
Saliency predictions



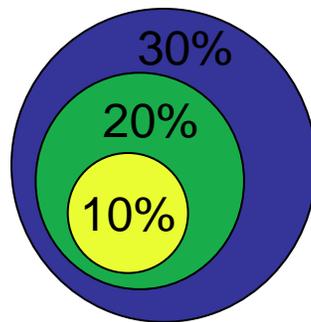
Comparison regions of interest



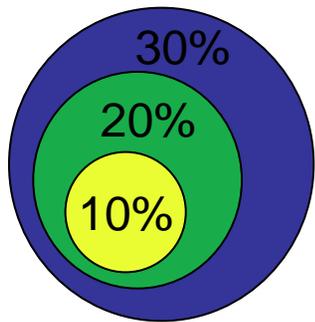
Saliency predictions



Saliency and Global scene priors



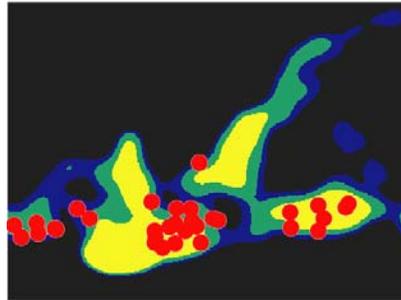
Comparison regions of interest



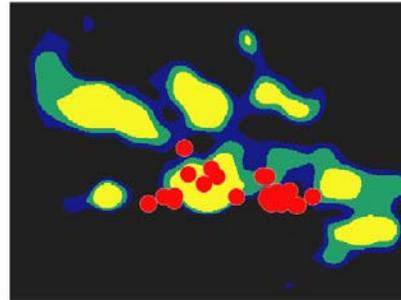
Saliency
predictions



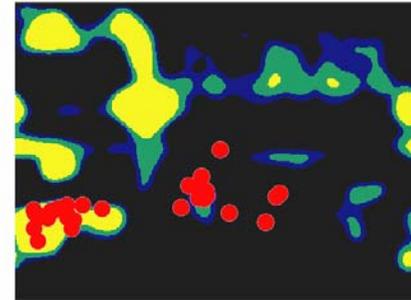
78%



25%

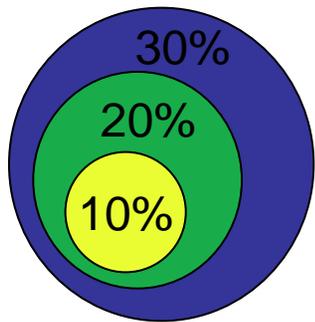


63%



Dots correspond to fixations 1-4

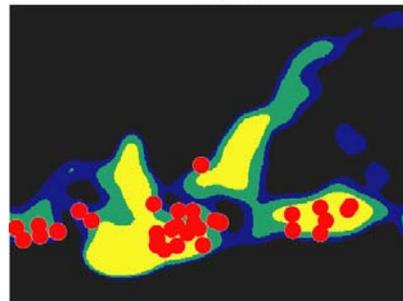
Comparison regions of interest



Saliency predictions



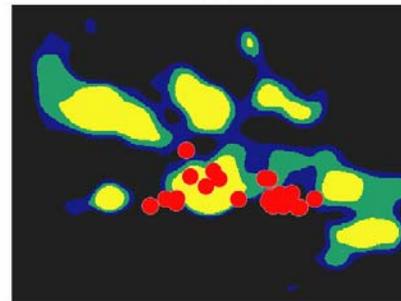
78%



96%



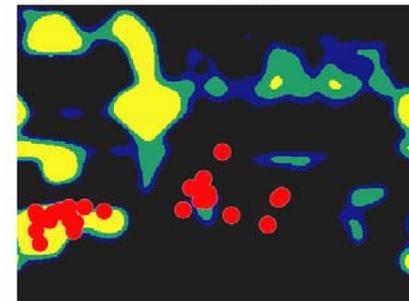
25%



96%

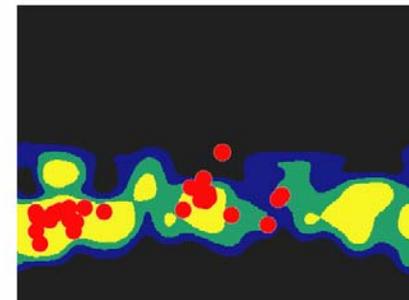
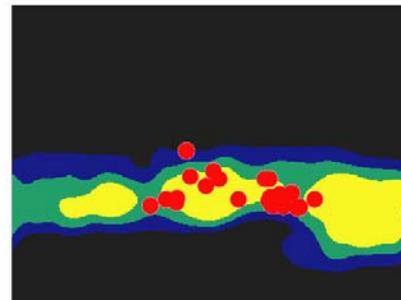
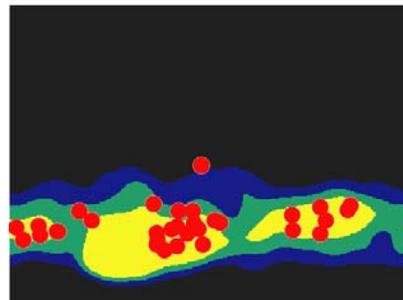


63%



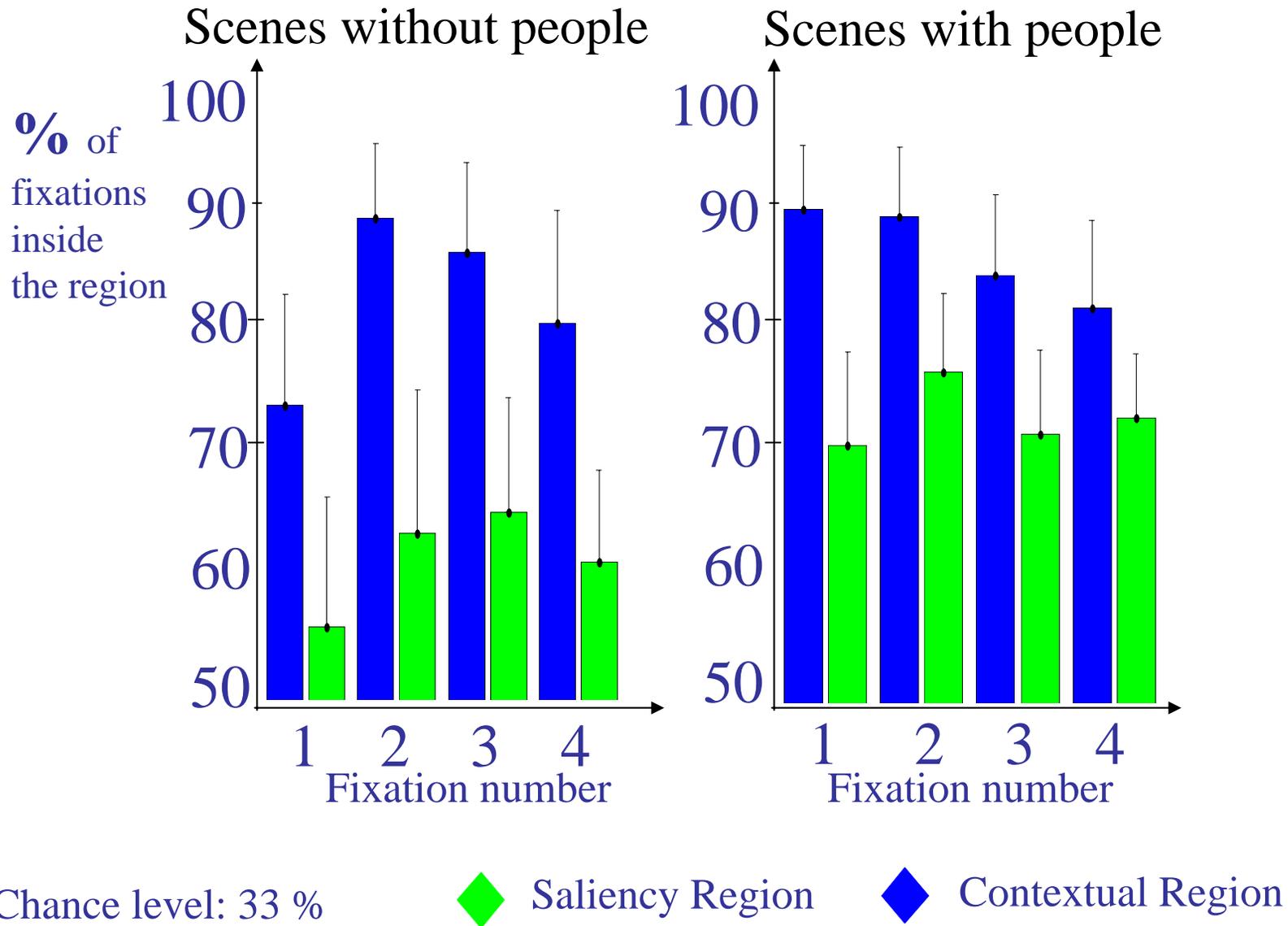
90%

Saliency and
Global scene
priors

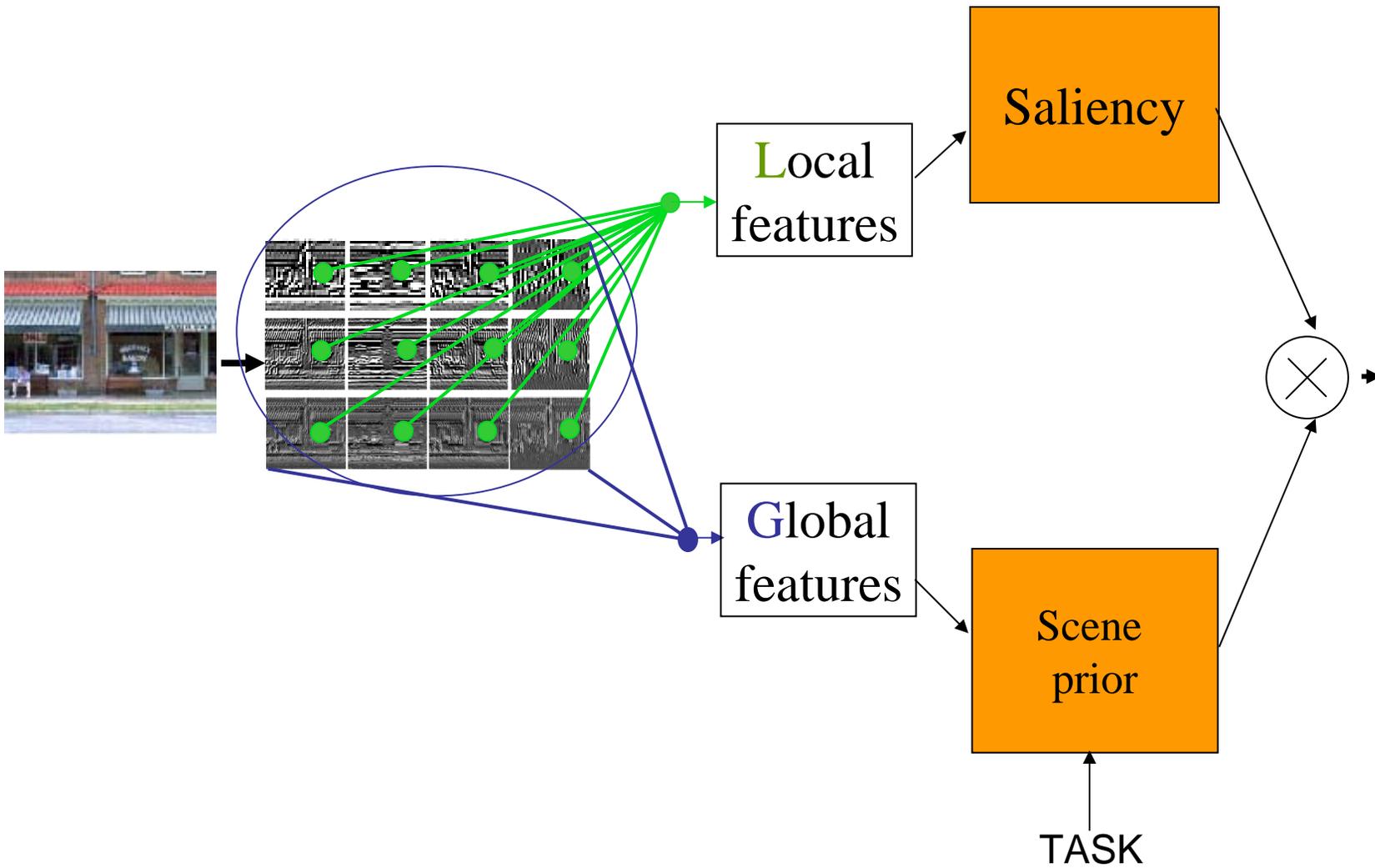


Dots correspond to fixations 1-4

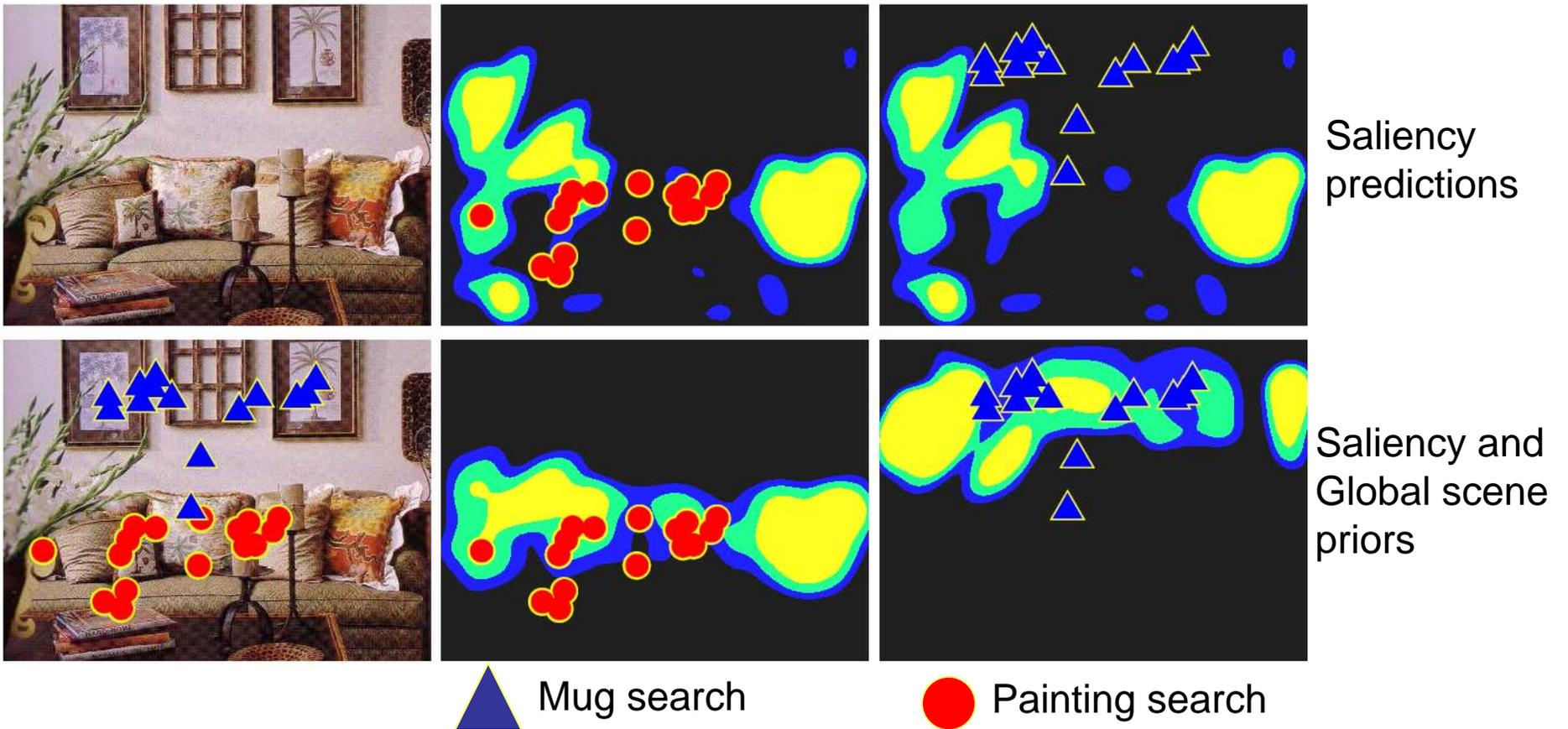
Results



Task modulation



Task modulation



Discussion

- From the computational perspective, scene context can be derived from global image properties and predict where objects are most likely to be.
- Scene context considerably improves predictions of fixation locations. A complete model of attention guidance in natural scenes requires both saliency and contextual pathways