

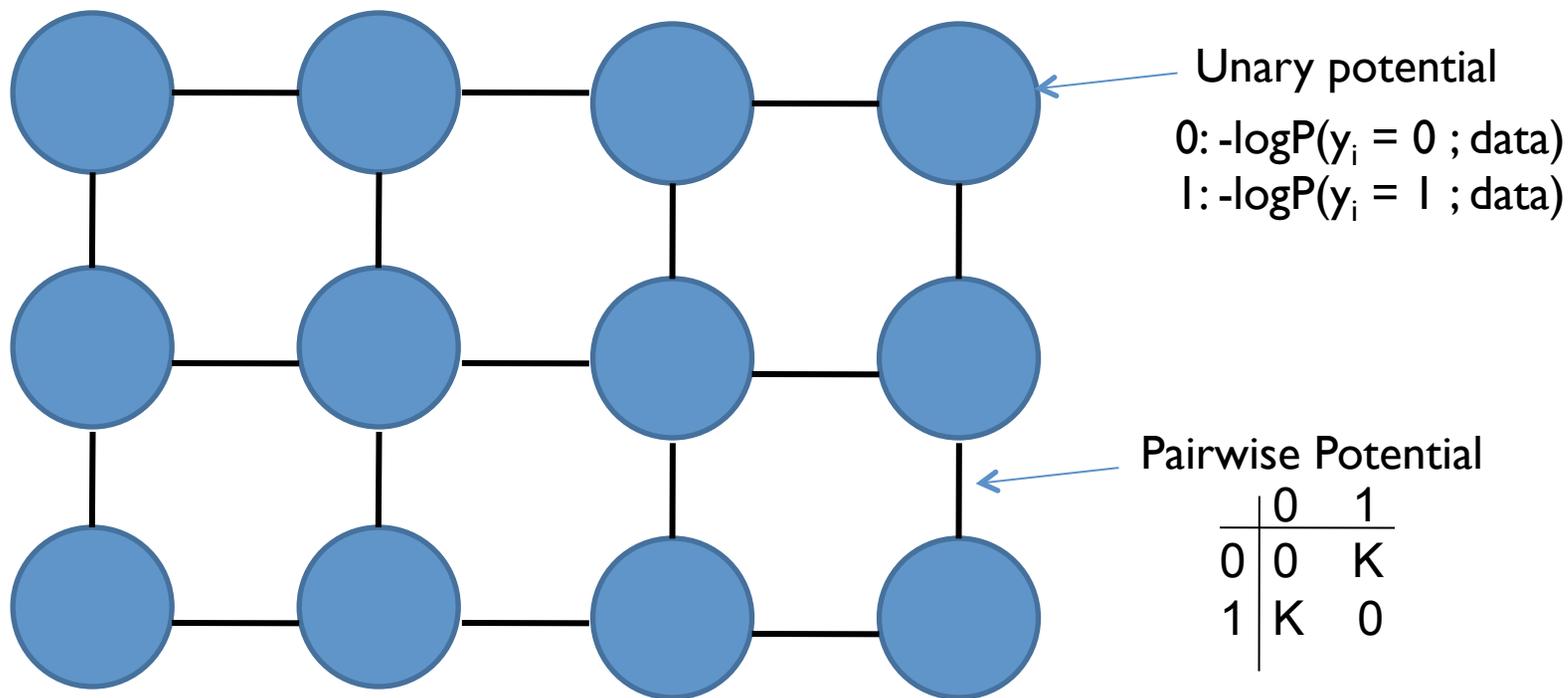
# **BIL 717**

# **Image Processing**

Erkut Erdem  
Dept. of Computer Engineering  
Hacettepe University

## **Semantic Segmentation**

# Review - Markov Random Fields



- Example: “label smoothing” grid

$$Energy(\mathbf{y}; \theta, data) = \sum_i \psi_1(y_i; \theta, data) + \sum_{i,j \in edges} \psi_2(y_i, y_j; \theta, data)$$

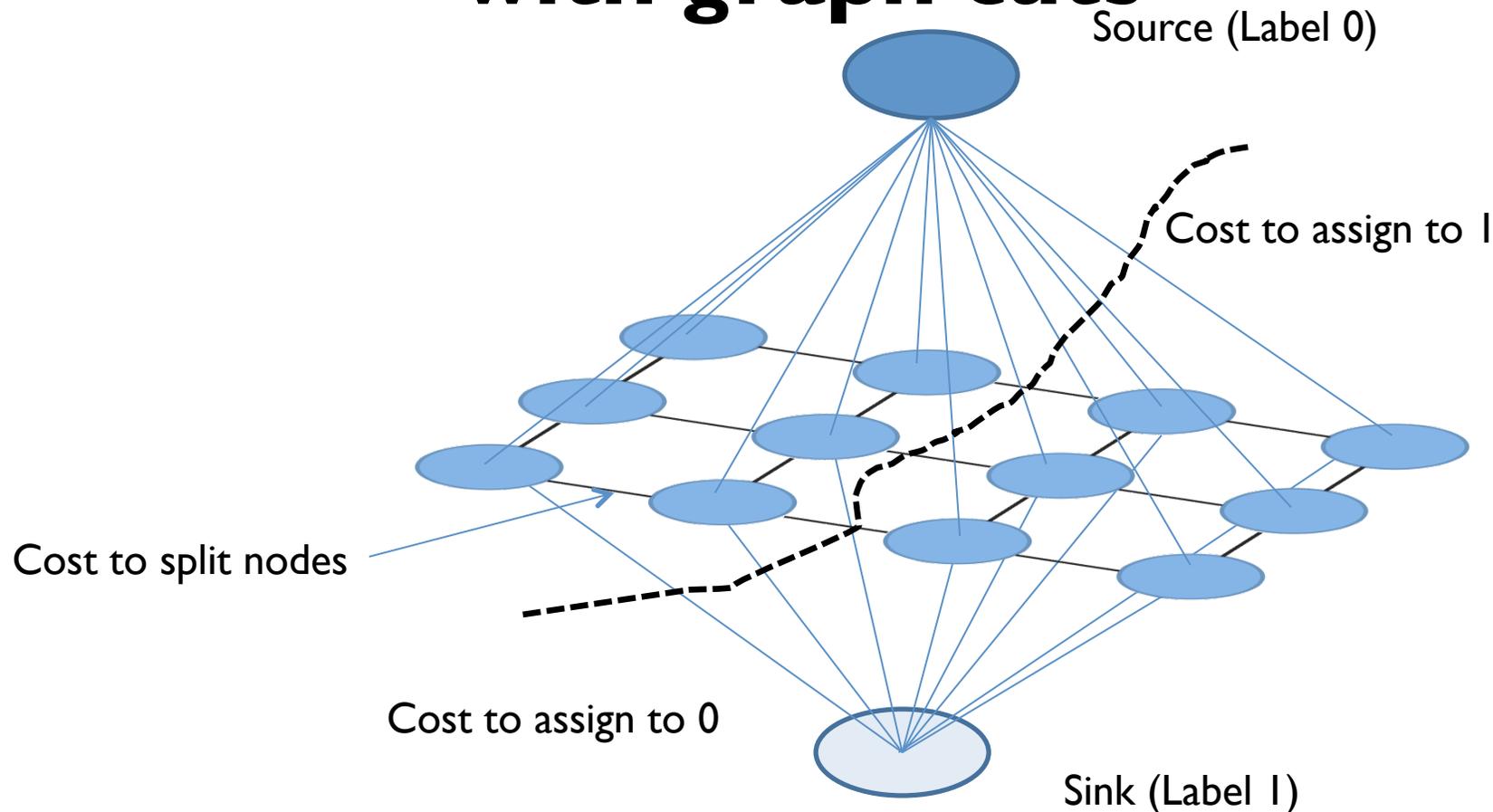
# Review - Solving MRFs with graph cuts

## Main idea:

- Construct a graph such that every *st*-cut corresponds to a joint assignment to the variables  $\mathbf{y}$
- The cost of the cut should be equal to the energy of the assignment,  $E(\mathbf{y}; \text{data})^*$ .
- The minimum-cut then corresponds to the minimum energy assignment,  $\mathbf{y}^* = \operatorname{argmin}_{\mathbf{y}} E(\mathbf{y}; \text{data})$ .

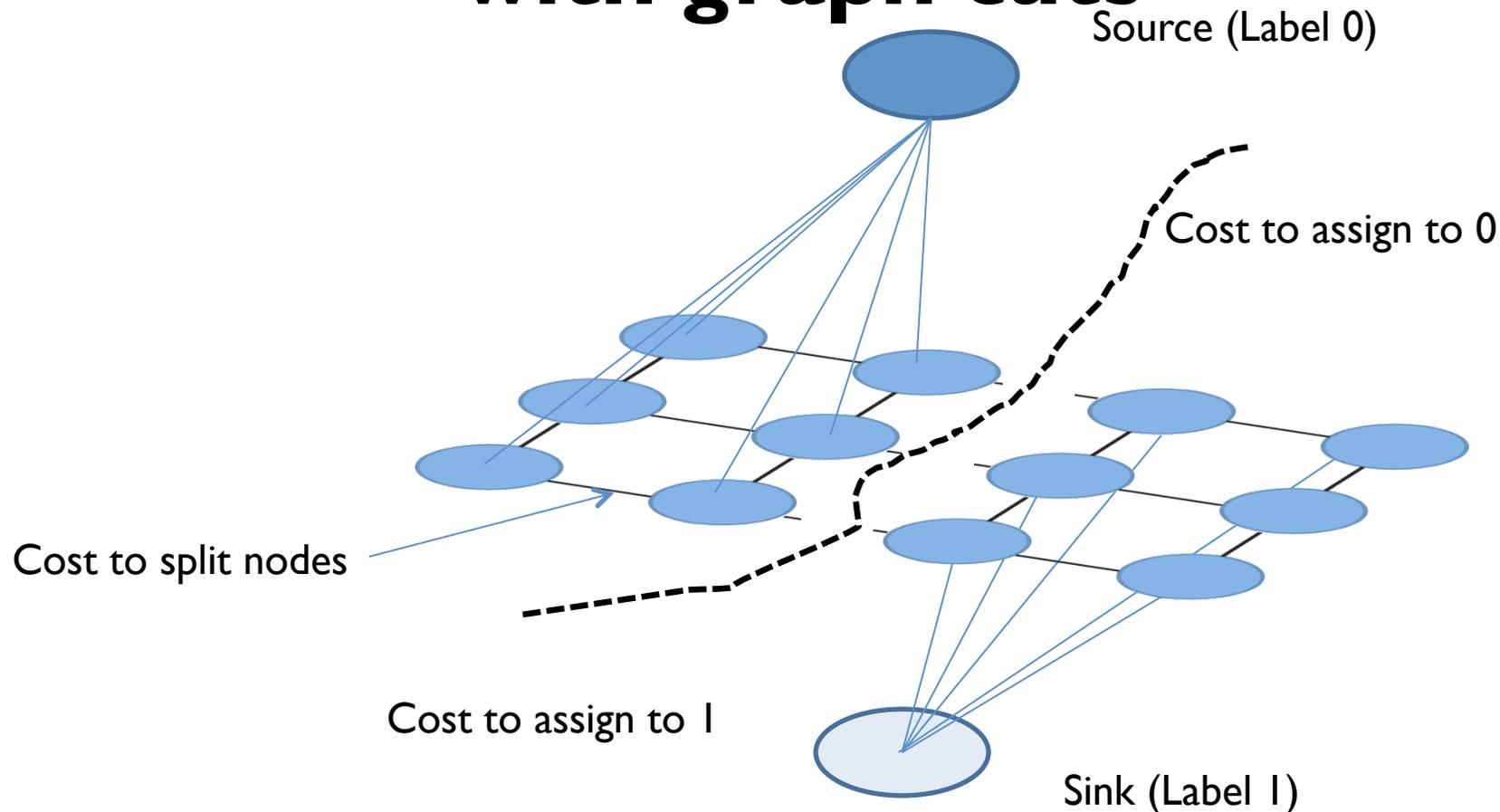
\* Requires non-negative energies

# Review - Solving MRFs with graph cuts



$$Energy(\mathbf{y}; \theta, data) = \sum_i \psi_1(y_i; \theta, data) + \sum_{i,j \in edges} \psi_2(y_i, y_j; \theta, data)$$

# Review - Solving MRFs with graph cuts



$$Energy(\mathbf{y}; \theta, data) = \sum_i \psi_1(y_i; \theta, data) + \sum_{i,j \in edges} \psi_2(y_i, y_j; \theta, data)$$

# Code for Image Segmentation

$$E(x) = \sum_i c_i x_i + \sum_{i,j} d_{ij} |x_i - x_j|$$

$$\begin{aligned} E: \{0,1\}^n &\rightarrow \mathbb{R} \\ 0 &\rightarrow \text{fg} \\ 1 &\rightarrow \text{bg} \end{aligned}$$

$n$  = number of pixels



$$x^* = \arg \min_x E(x)$$

How to minimize  $E(x)$ ?

Global Minimum ( $x^*$ )

# Review - How does the code look like?

```
Graph *g;
```

```
For all pixels p
```

```
    /* Add a node to the graph */  
    nodeID(p) = g->add_node();
```

```
    /* Set cost of terminal edges */  
    set_weights(nodeID(p), fgCost(p),  
                bgCost(p));
```

```
end
```

```
for all adjacent pixels p,q  
    add_weights(nodeID(p), nodeID(q),  
                cost(p,q));
```

```
end
```

```
g->compute_maxflow();
```

```
label_p = g->is_connected_to_source(nodeID(p));  
// is the label of pixel p (0 or 1)
```

 Source (0)

 Sink (1)

# Review - How does the code look like?

```
Graph *g;
```

```
For all pixels p
```

```
/* Add a node to the graph */  
nodeID(p) = g->add_node();
```

```
/* Set cost of terminal edges */  
set_weights(nodeID(p), fgCost(p),  
            bgCost(p));
```

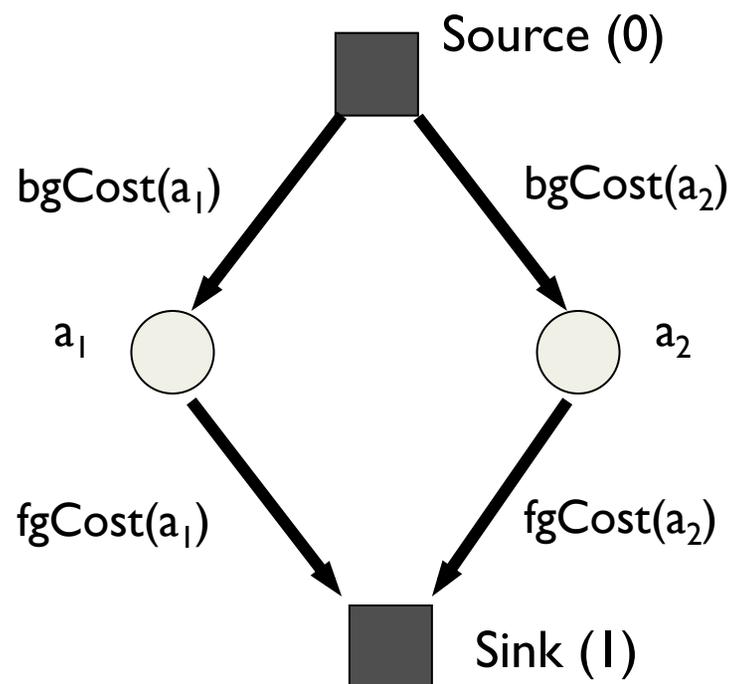
```
end
```

```
for all adjacent pixels p,q  
    add_weights(nodeID(p), nodeID(q),  
                cost(p,q));
```

```
end
```

```
g->compute_maxflow();
```

```
label_p = g->is_connected_to_source(nodeID(p));  
// is the label of pixel p (0 or 1)
```



# Review - How does the code look like?

```
Graph *g;
```

```
For all pixels p
```

```
/* Add a node to the graph */  
nodeID(p) = g->add_node();
```

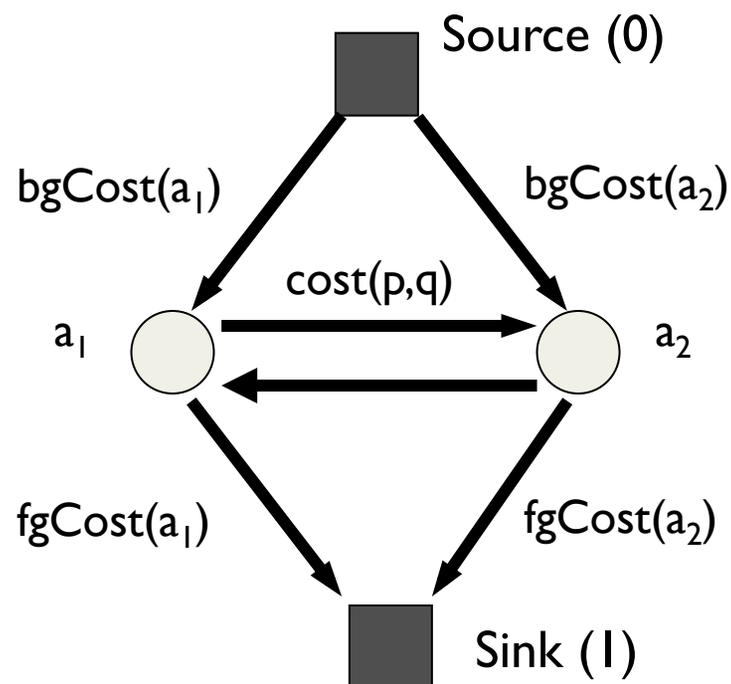
```
/* Set cost of terminal edges */  
set_weights(nodeID(p), fgCost(p),  
            bgCost(p));
```

```
end
```

```
for all adjacent pixels p,q  
    add_weights(nodeID(p), nodeID(q),  
                cost(p,q));  
end
```

```
g->compute_maxflow();
```

```
label_p = g->is_connected_to_source(nodeID(p));  
// is the label of pixel p (0 or 1)
```



# Review - How does the code look like?

```
Graph *g;
```

```
For all pixels p
```

```
/* Add a node to the graph */  
nodeID(p) = g->add_node();
```

```
/* Set cost of terminal edges */  
set_weights(nodeID(p), fgCost(p),  
            bgCost(p));
```

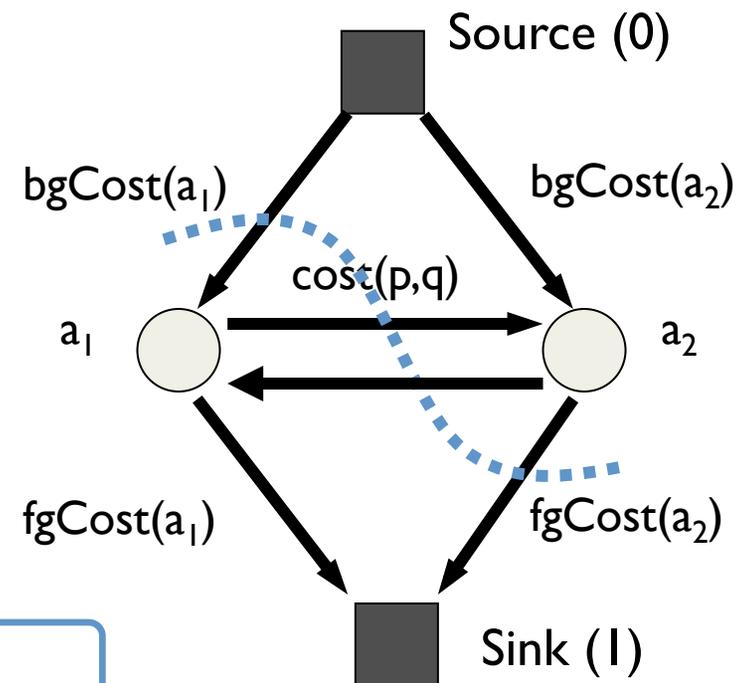
```
end
```

```
for all adjacent pixels p,q  
    add_weights(nodeID(p), nodeID(q),  
               cost(p,q));
```

```
end
```

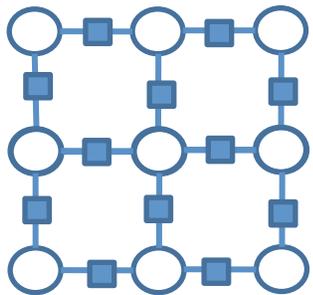
```
g->compute_maxflow();
```

```
label_p = g->is_connected_to_source(nodeID(p));  
// is the label of pixel p (0 or 1)
```



$$a_1 = bg \quad a_2 = fg$$

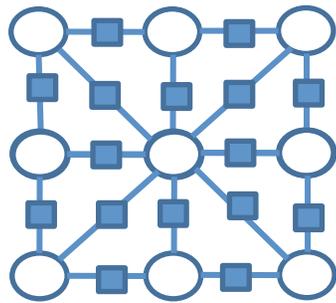
# Review - Random Fields in Vision



4-connected;  
pairwise MRF

$$E(\mathbf{x}) = \sum_{i,j \in N_4} \theta_{ij}(\mathbf{x}_i, \mathbf{x}_j)$$

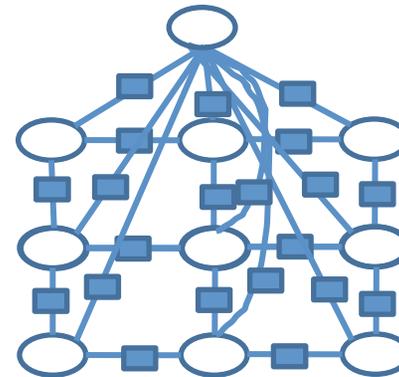
Order 2



higher(8)-connected;  
pairwise MRF

$$E(\mathbf{x}) = \sum_{i,j \in N_8} \theta_{ij}(\mathbf{x}_i, \mathbf{x}_j)$$

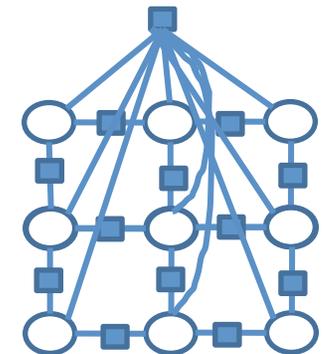
Order 2



MRF with  
global variables

$$E(\mathbf{x}) = \sum_{i,j \in N_8} \theta_{ij}(\mathbf{x}_i, \mathbf{x}_j)$$

Order 2



Higher-order MRF

$$E(\mathbf{x}) = \sum_{i,j \in N_4} \theta_{ij}(\mathbf{x}_i, \mathbf{x}_j) + \theta(\mathbf{x}_1, \dots, \mathbf{x}_n)$$

Order n

# Review - MRF with global potential

GrabCut model [Rother et. al. '04]

$$E(x, \theta^F, \theta^B) = \sum_i F_i(\theta^F)x_i + B_i(\theta^B)(1-x_i) + \sum_{i,j \in N} |x_i - x_j|$$

$$F_i = -\log \Pr(z_i | \theta^F)$$

$$B_i = -\log \Pr(z_i | \theta^B)$$

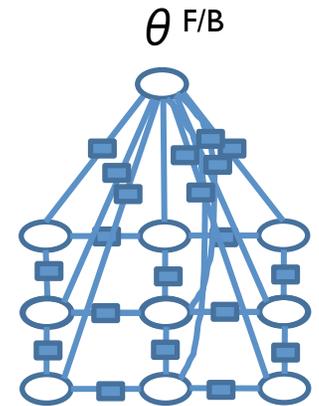
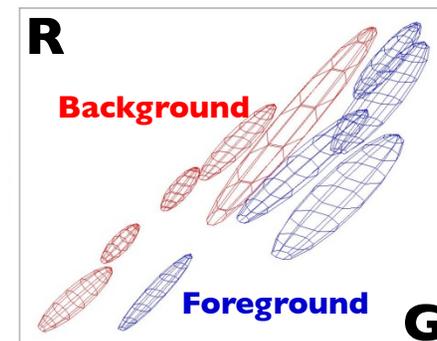


Image z



Output x

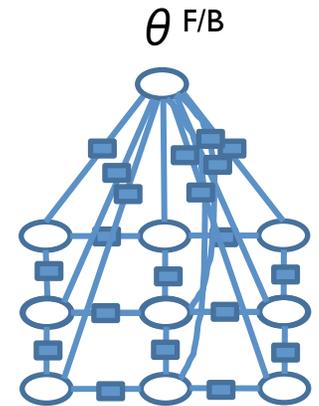


$\theta^{F/B}$  Gaussian Mixture models

**Problem:** for unknown  $x, \theta^F, \theta^B$  the optimization is NP-hard! [Vicente et al. '09]

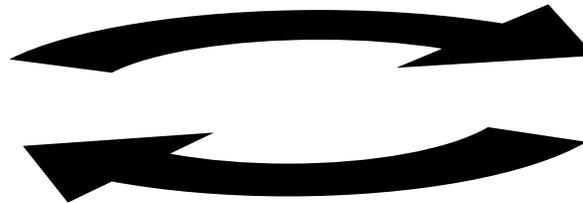
# Review - GrabCut: Iterated Graph Cuts

[Rother et al. Siggraph '04]



$$\min_{\theta^F, \theta^B} E(x, \theta^F, \theta^B)$$

**Learning of the  
colour distributions**

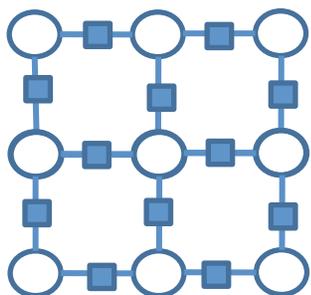


$$\min_x E(x, \theta^F, \theta^B)$$

**Graph cut to infer  
segmentation**

Most systems with global variables work like that  
e.g. [ObjCut Kumar et al. '05, PoseCut Bray et al. '06, LayoutCRF Winn et al. '06]

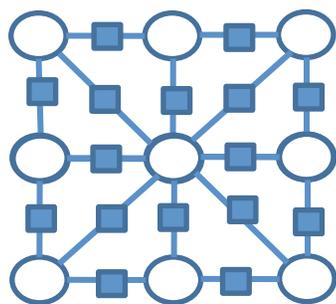
# Review - Random Fields in Vision



4-connected;  
pairwise MRF

$$E(x) = \sum_{i,j \in N_4} \theta_{ij}(x_i, x_j)$$

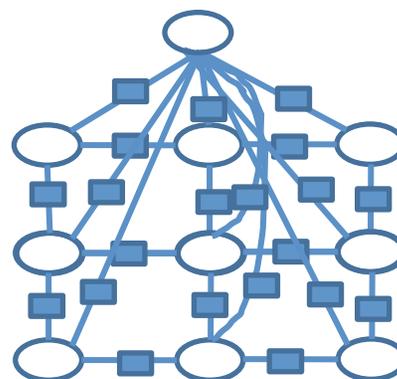
Order 2



higher(8)-connected;  
pairwise MRF

$$E(x) = \sum_{i,j \in N_8} \theta_{ij}(x_i, x_j)$$

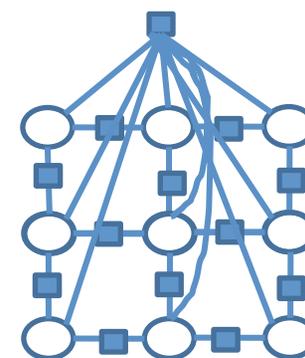
Order 2



MRF with  
global variables

$$E(x) = \sum_{i,j \in N_8} \theta_{ij}(x_i, x_j)$$

Order 2



Higher-order MRF

$$E(x) = \sum_{i,j \in N_4} \theta_{ij}(x_i, x_j) + \theta(x_1, \dots, x_n)$$

Order n

# Review - Why Higher-order Functions?

In general  $\theta(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \neq \theta(\mathbf{x}_1, \mathbf{x}_2) + \theta(\mathbf{x}_1, \mathbf{x}_3) + \theta(\mathbf{x}_2, \mathbf{x}_3)$

Reasons for higher-order RFs:

1. Even better image(texture) models:

- Field-of Expert [FoE, Roth et al. '05]
- Curvature [Woodford et al. '08]

2. Use **global** Priors:

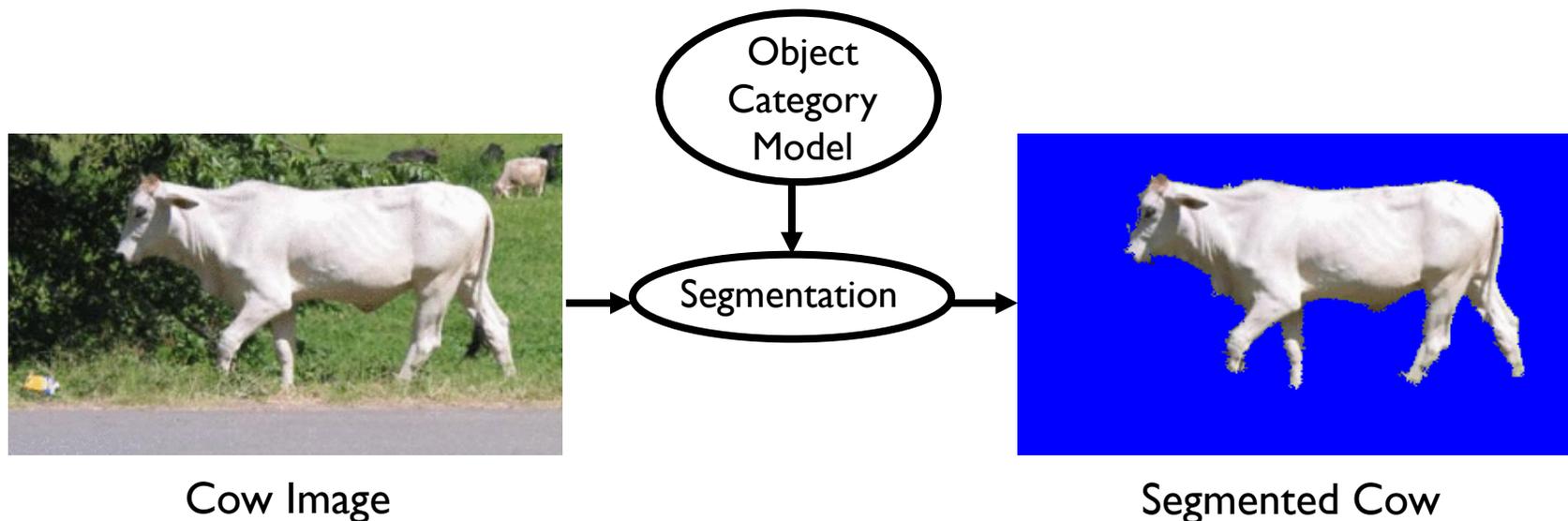
- Connectivity [Vicente et al. '08, Nowozin et al. '09]
- Better encoding label statistics [Woodford et al. '09]
- Convert global variables to global factors [Vicente et al. '09]

# Semantic Segmentation

- Joint recognition & segmentation
  - segmenting all the objects in a given image and identifying their visual categories
- aka scene parsing or image parsing
- Early studies aim at segmenting out a single object of a known category
  - Borenstein & Ullman, 2002, Liebe & Schiele, 2003,

# Early Studies of Semantic Segmentation

- Given an image and object category, to segment the object



- Segmentation should (ideally) be
  - shaped like the object e.g. cow-like
  - obtained efficiently in an unsupervised manner
  - able to handle self-occlusion

# Early Studies of Semantic Segmentation



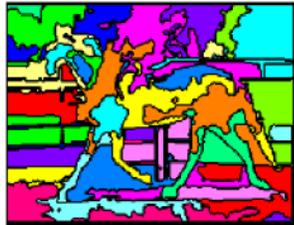
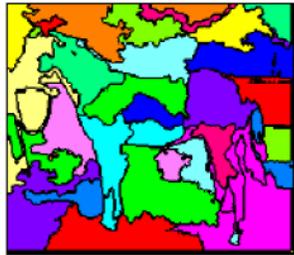
# Early Studies of Semantic Segmentation



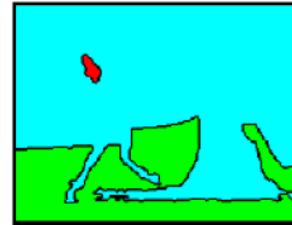
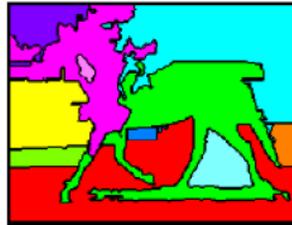
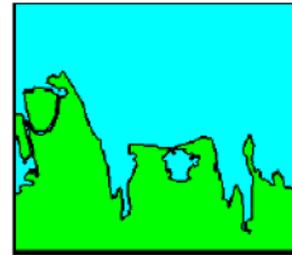
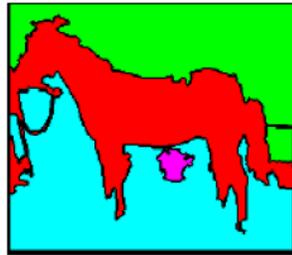
# Early Studies of Semantic Segmentation

Using Normalized Cuts, Shi & Malik, 1997

Input



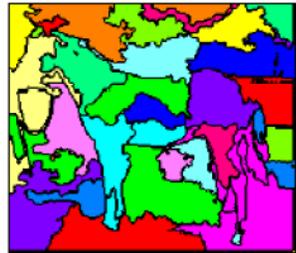
Bottom-up



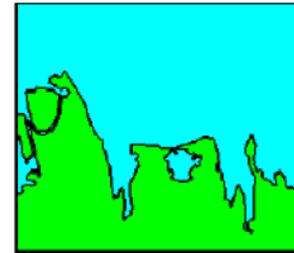
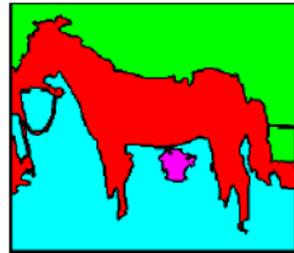
# Early Studies of Semantic Segmentation

Using Normalized Cuts, Shi & Malik, 1997

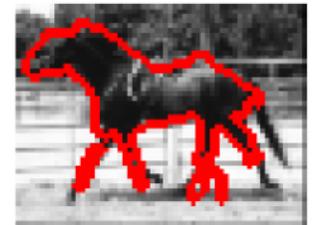
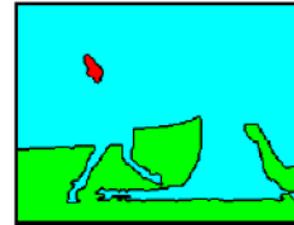
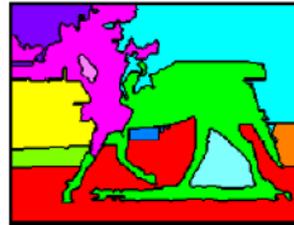
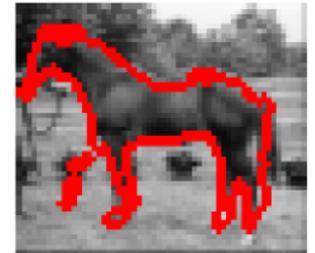
Input



Bottom-up

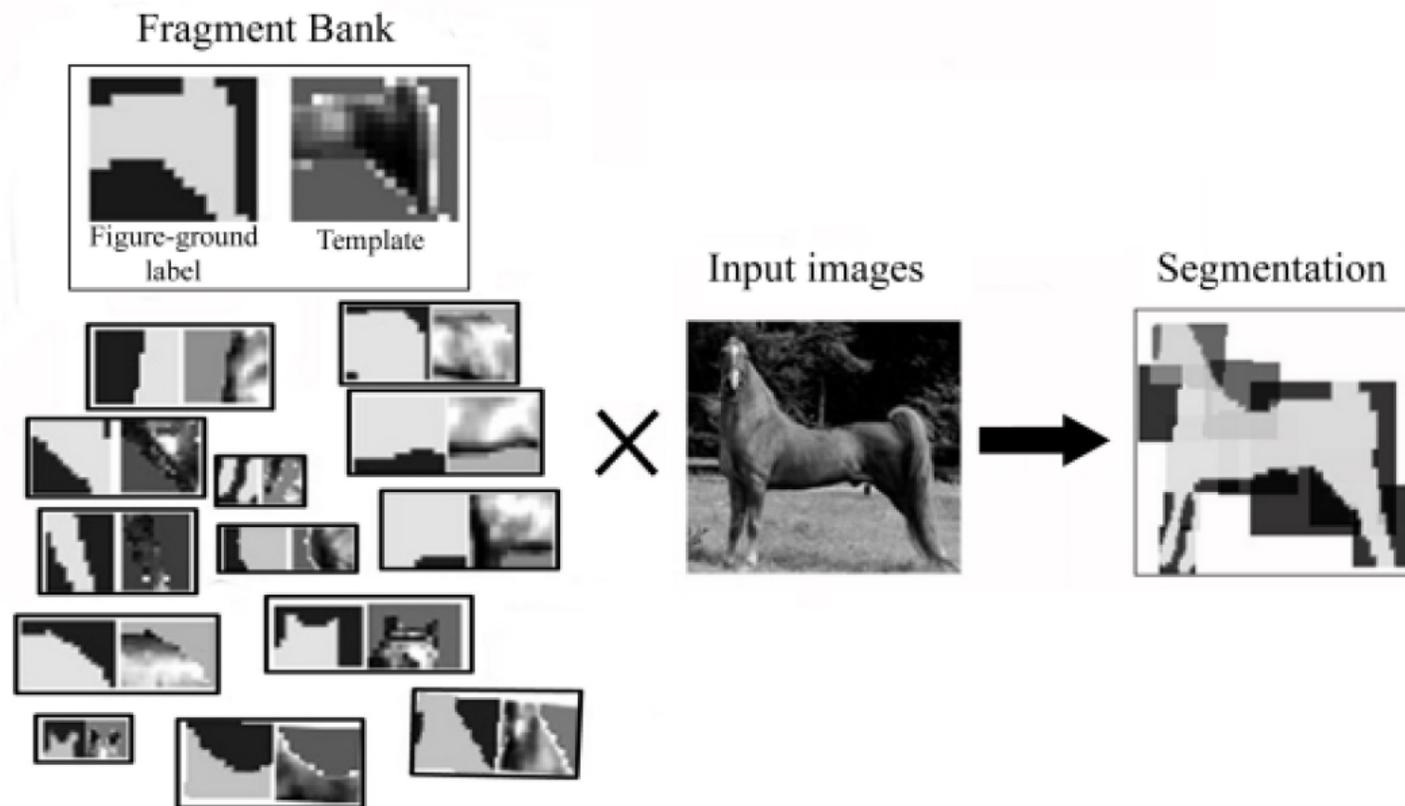


Top-down

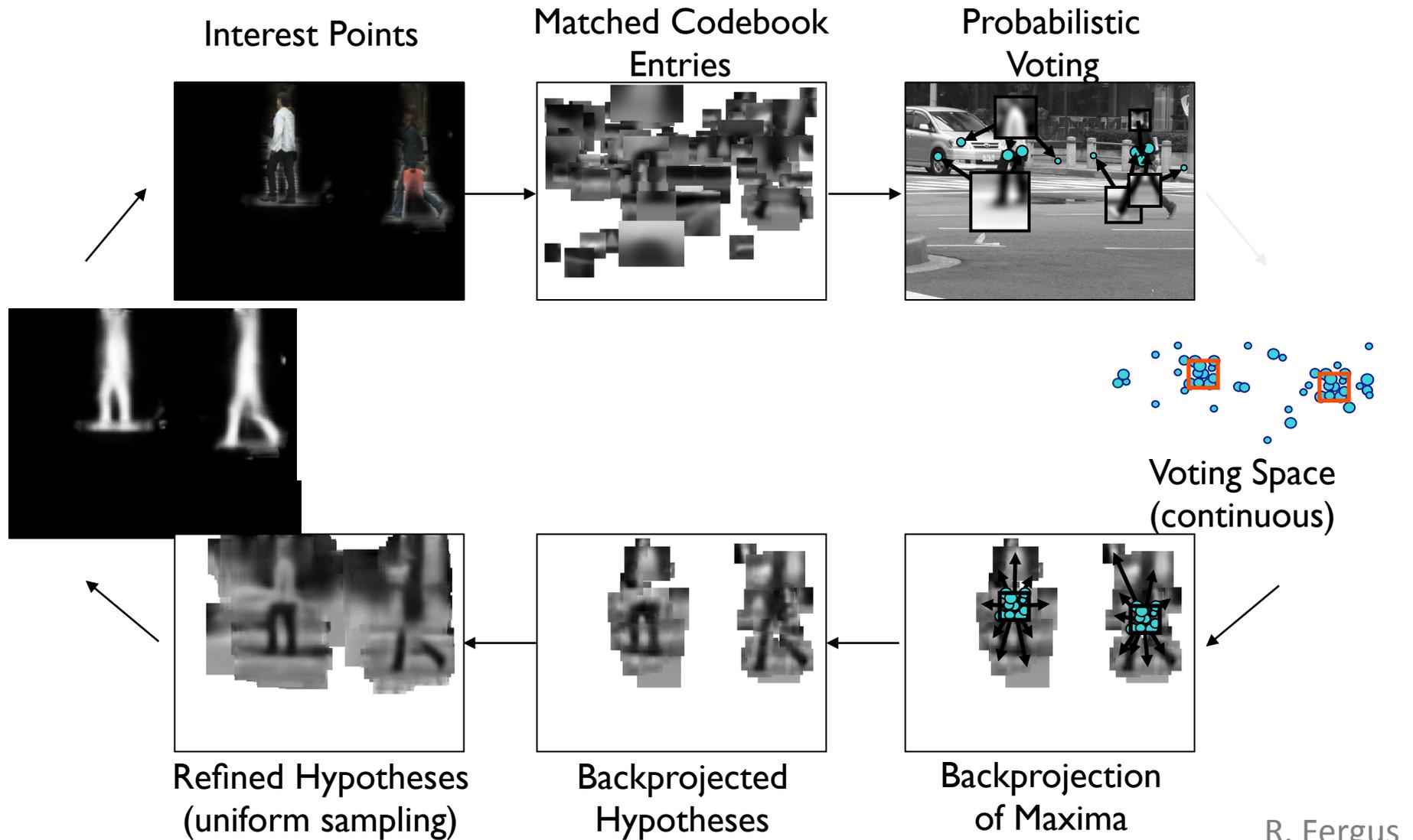


Borenstein and Ullman, ECCV 2002

# Jigsaw approach: Borenstein and Ullman, 2002



# Implicit Shape Model - Liebe and Schiele, 2003



# Random Fields for segmentation

$I$  = Image pixels (observed)

$h$  = foreground/background labels (hidden) – one label per pixel

$\theta$  = Parameters

$$\underbrace{p(h | I, \theta)}$$

**Posterior**

# Random Fields for segmentation

$I$  = Image pixels (observed)

$h$  = foreground/background labels (hidden) – one label per pixel

$\theta$  = Parameters

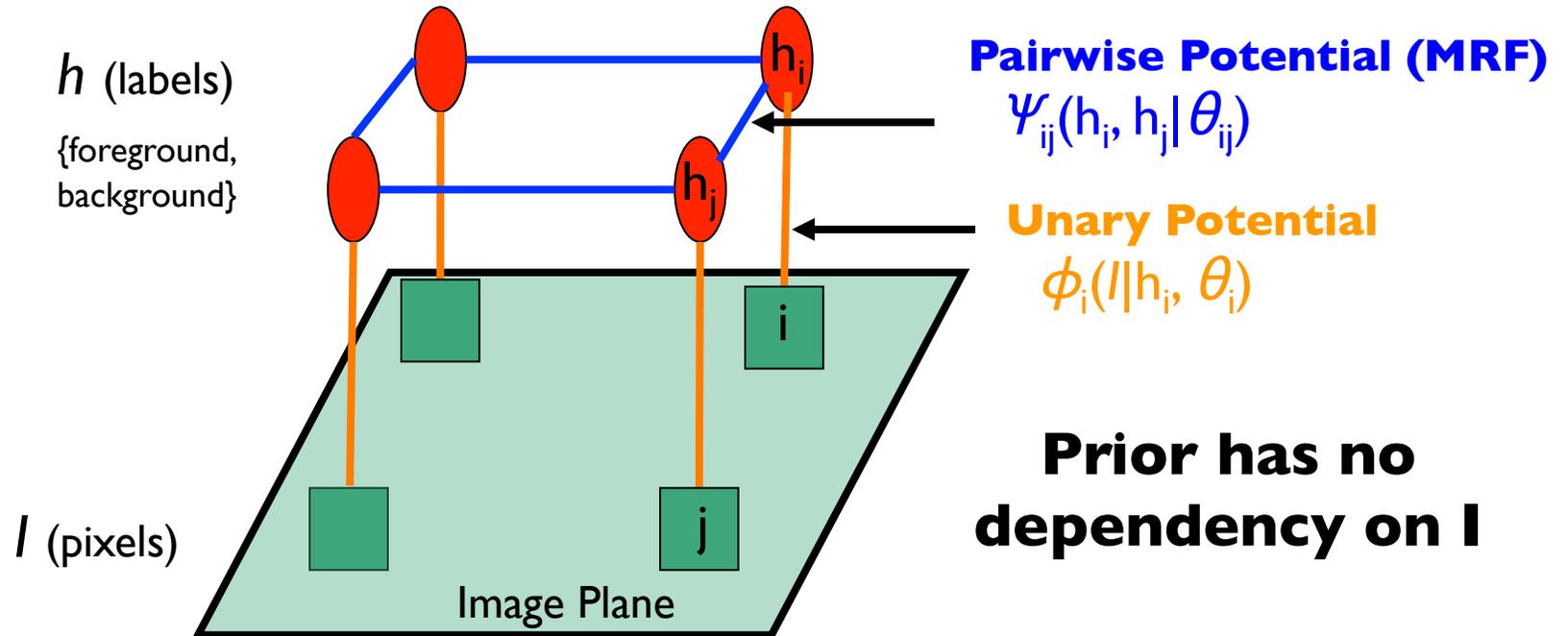
$$\underbrace{p(h | I, \theta)}_{\text{Posterior}} \propto \underbrace{p(I, h | \theta)}_{\text{Joint}} = \underbrace{p(I | h, \theta)}_{\text{Likelihood}} \underbrace{p(h | \theta)}_{\text{Prior}}$$

1. Generative approach models joint  
→ Markov random field (MRF)
2. Discriminative approach models posterior directly  
→ Conditional random field (CRF)

# Generative Markov Random Field

$$p(h, I | \theta) = p(I | h, \theta) p(h | \theta)$$

$$= \frac{1}{Z(\theta)} \left[ \underbrace{\prod_i \phi_i(I | h_i, \theta_i)}_{\text{Likelihood}} \underbrace{\prod_{ij} \psi_{ij}(h_i, h_j | \theta_{ij})}_{\text{MRF Prior}} \right]$$



# Conditional Random Field

Lafferty, McCallum and Pereira 2001

Discriminative approach

$$p(h | I, \theta) = \frac{1}{Z(I, \theta)} \left[ \underbrace{\prod_i \phi_i(h_i, I | \theta_i)}_{\text{Unary}} \underbrace{\prod_{ij} \psi_{ij}(h_i, h_j, I | \theta_{ij})}_{\text{Pairwise}} \right]$$

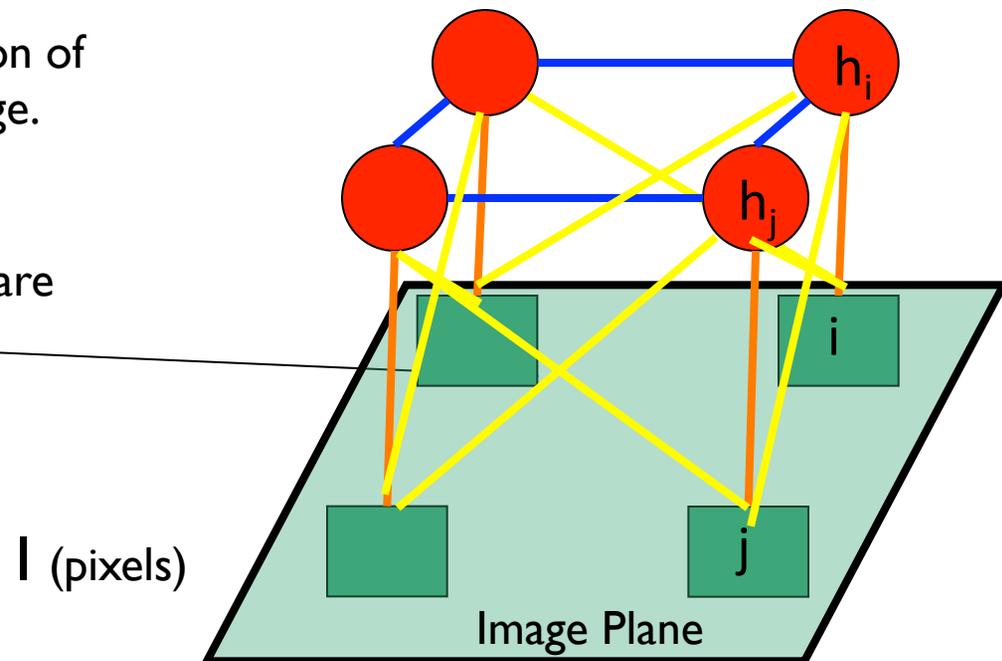
Unary

Pairwise

- Dependency on  $I$  allows introduction of pairwise terms that make use of image.

- For example, neighboring labels should be similar only if pixel colors are similar  $\rightarrow$  Contrast term

e.g Kumar and Hebert 2003

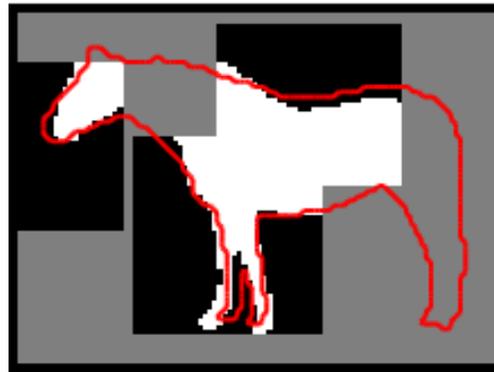
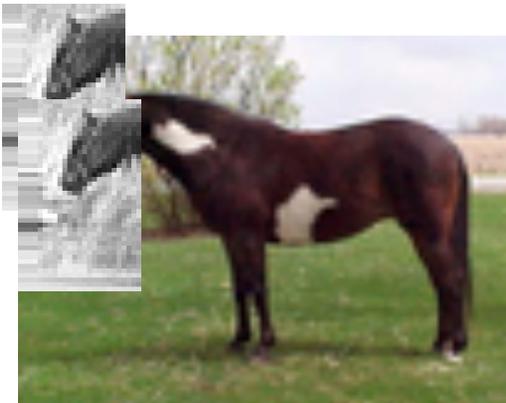


# Levin & Weiss [ECCV 2006]

$$E(h; I) = \sum_i \lambda_i |h - h_{F_i, I}| + \sum_{ij} w(i, j) |h_i - h_j|$$

Consistency with  
fragments  
segmentation

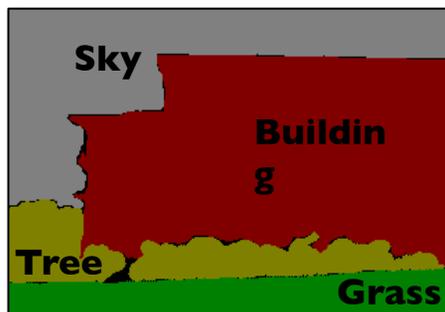
Segmentation  
alignment with  
image edges



Resulting min-cut  
segmentation

# Semantic Segmentation

## Joint Object recognition & segmentation

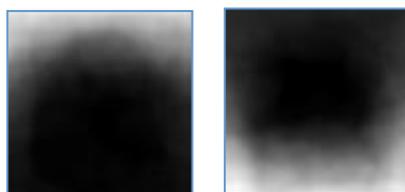


$$E(x, \omega) = \sum_i \theta_i(\omega, x_i) + \sum_i \theta_i(x_i) + \sum_i \theta_i(x_i) + \sum_{i,j} \theta_{ij}(x_i, x_j)$$

(color)
(location)
(class)
(edge aware Ising prior)

$x_i \in \{1, \dots, K\}$  for  $K$  object classes

Location



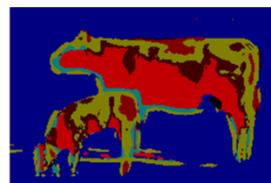
sky

grass

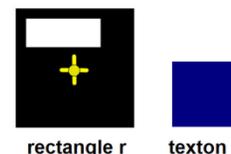
Class (boosted textons)



(a) Input image



(b) Texton map



(c) Feature pair = (r,t)



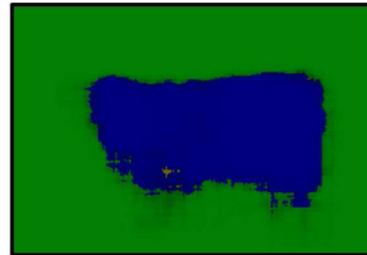
(d) Superimposed rectangles

# Semantic Segmentation

## Joint Object recognition & segmentation

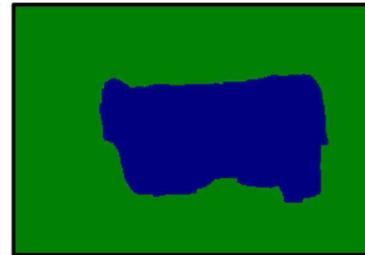


(a)



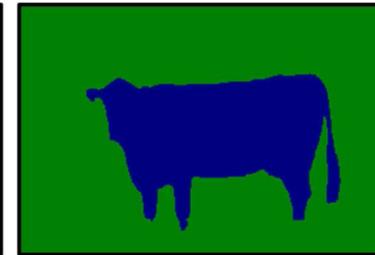
(b) 69.6%

Class+  
location



(c) 70.3%

+ edges



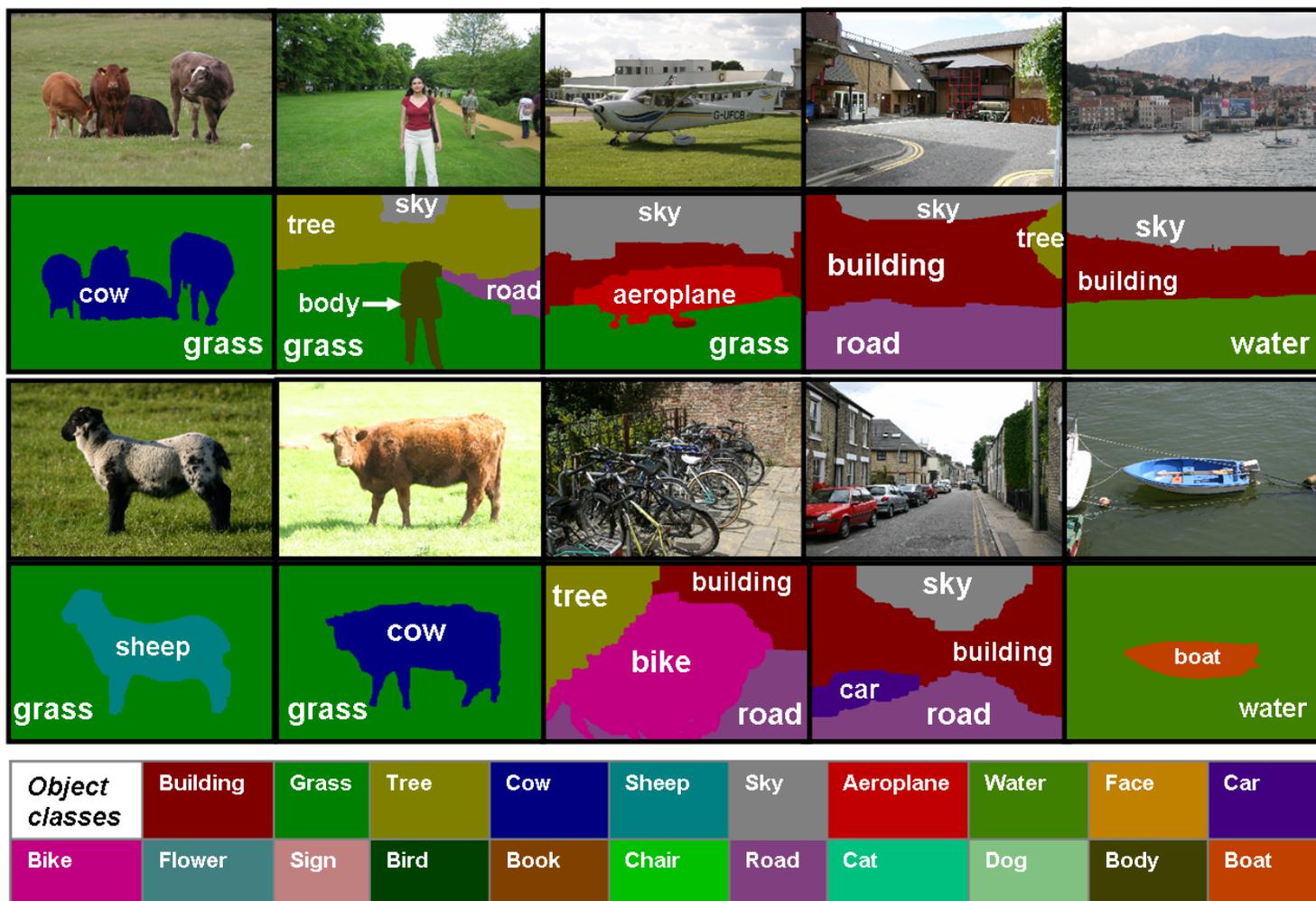
(d) 72.2%

+ color

# Semantic Segmentation

## Joint Object recognition & segmentation

Good results ...



[TextonBoost; Shotton et al, '06]

C. Rother

# Semantic Segmentation

## Joint Object recognition & segmentation

Failure cases...



# Semantic Segmentation

## Joint Object recognition & segmentation

Goal: Detect and segment test image:



Large set of example segmentation:



T(1)

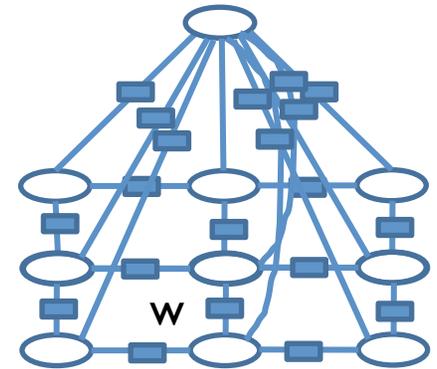


T(2)



T(3)

Up to 2.000.000 shape templates



$$E(x,w): \{0, 1\}^n \times \{\text{Exemplar}\} \rightarrow \mathbb{R}$$

$$E(x,w) = \sum_i |T(w)_i - x_i| + \sum_{i,j \in \mathbf{N}_4} \theta_{ij}(x_i, x_j)$$

“Hamming distance”

# Semantic Segmentation

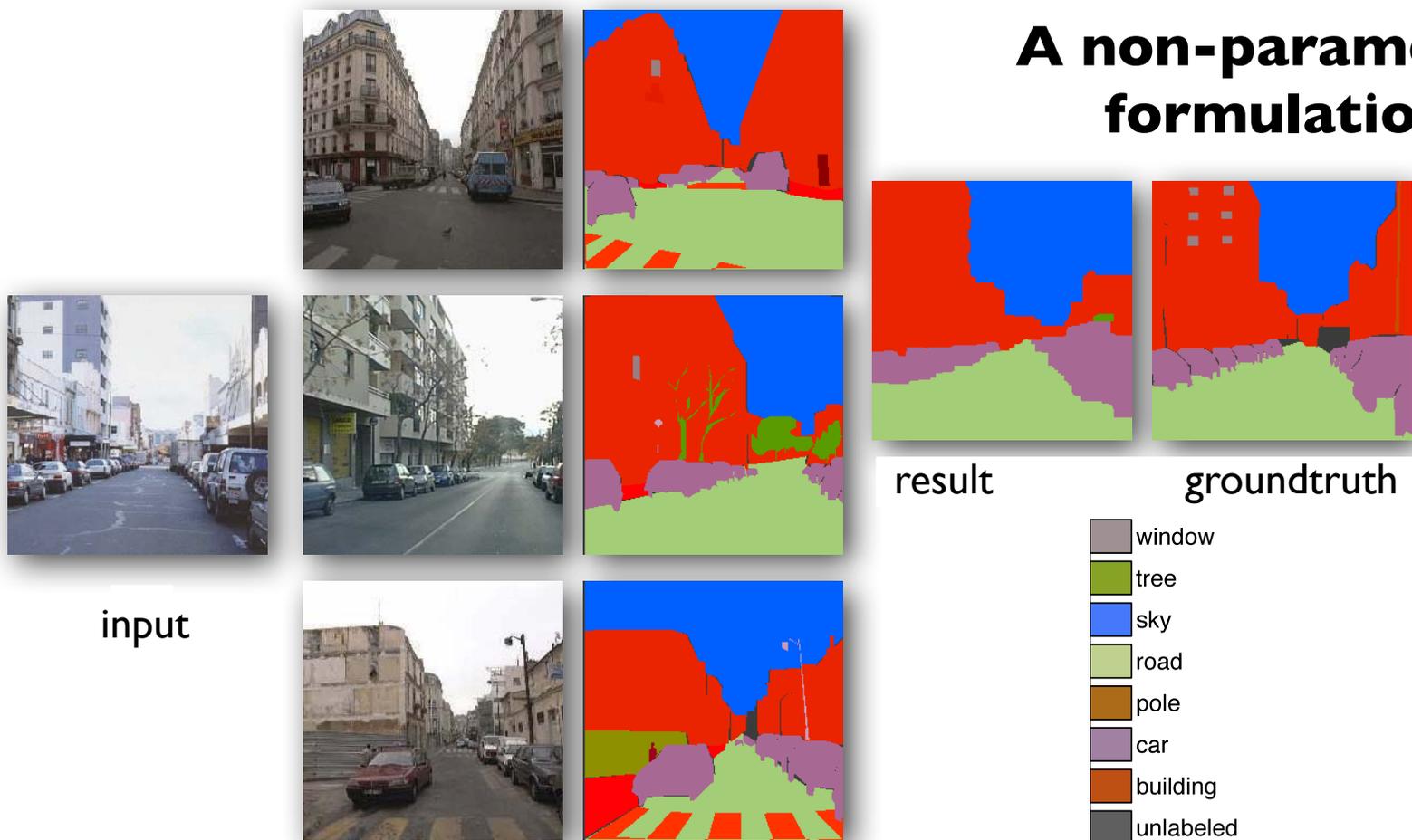
## Joint Object recognition & segmentation



UIUC dataset; 98.8% accuracy

# Nonparametric Scene Parsing via Label Transfer (Liu et al. TPAMI'12)

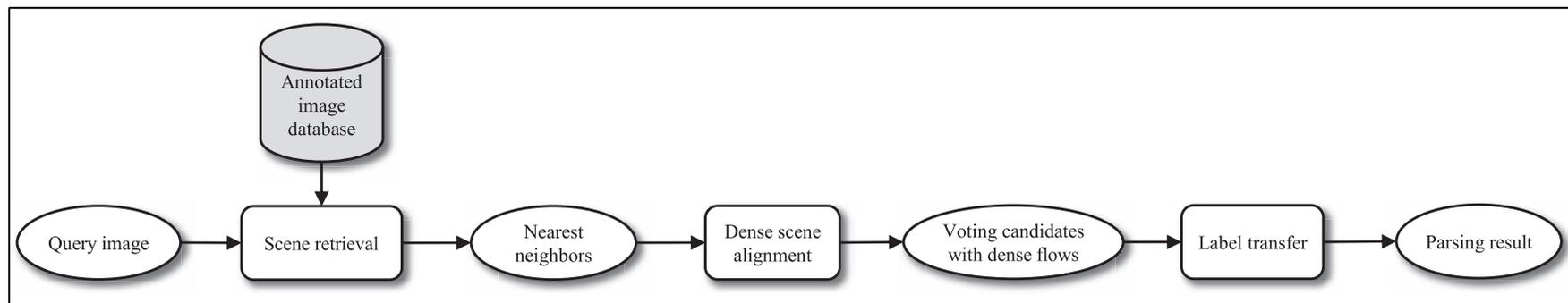
**A non-parametric formulation**



retrieved images and their annotations

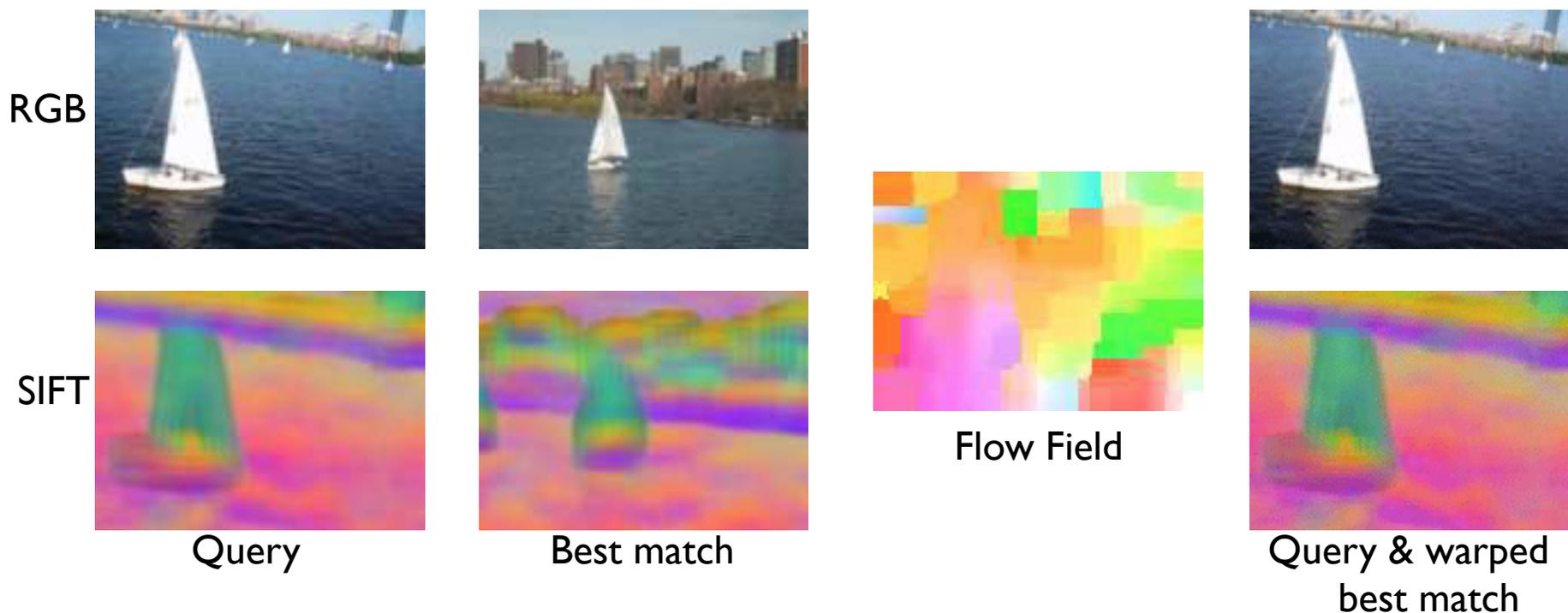
# Nonparametric Scene Parsing via Label Transfer

- Framework consists of three main modules:
  1. Scene retrieval: finding nearest neighbors (k-NN approach)
  2. Dense scene alignment: dense scene matching (SIFT Flow)
  3. Label transfer: using a MRF model to label input image



# Dense Scene Alignment via SIFT Flow

- SIFT Flow (Liu et al., ECCV 2008)
  - Finds semantically meaningful correspondences among two images by matching local SIFT descriptors



# Dense Scene Alignment via SIFT Flow

- SIFT Flow (Liu et al., ECCV 2008)
  - Finds semantically meaningful correspondences among two images by matching local SIFT descriptors

$$E(\mathbf{w}) = \sum_{\mathbf{p}} \min(\|s_1(\mathbf{p}) - s_2(\mathbf{p} + \mathbf{w}(\mathbf{p}))\|_1, t) + \quad \text{data term}$$

$$\sum_{\mathbf{p}} \eta(|u(\mathbf{p})| + |v(\mathbf{p})|) + \quad \text{small displacement term}$$

$$\sum_{(\mathbf{p}, \mathbf{q}) \in \mathcal{E}} \min(\lambda|u(\mathbf{p}) - u(\mathbf{q})|, d) + \quad \text{smoothness term}$$
$$\min(\lambda|v(\mathbf{p}) - v(\mathbf{q})|, d),$$

$\mathbf{w}(\mathbf{p}) = (u(\mathbf{p}), v(\mathbf{p}))$  : flow vector at point  $\mathbf{p}$

# Label Transfer

- A set of voting candidates  $\{s_i; c_i; w_i\}_{i=1:M}$  is obtained from the retrieved images with  $s_i$ ,  $c_i$ , and  $w_i$  denoting the SIFT image, annotation, and SIFT flow field of the  $i$ th voting candidate.
- A probabilistic MRF model is built to integrate
  - multiple category labels,
  - prior object (category) information
  - spatial smoothness of category labels

$$- \log P(c|I, s, \{s_i, c_i, w_i\}) = \sum_{\mathbf{p}} \psi(c(\mathbf{p}); s, \{s'_i\})$$
$$+ \alpha \sum_{\mathbf{p}} \lambda(c(\mathbf{p})) + \beta \sum_{\{\mathbf{p}, \mathbf{q}\} \in \mathcal{E}} \phi(c(\mathbf{p}), c(\mathbf{q}); I) + \log Z$$

# Label Transfer

- Likelihood term:

$$\psi(c(\mathbf{p}) = l) = \begin{cases} \min_{i \in \Omega_{\mathbf{p},l}} \|s(\mathbf{p}) - s_i(\mathbf{p} + \mathbf{w}(\mathbf{p}))\|, & \Omega_{\mathbf{p},l} \neq \emptyset, \\ \tau, & \Omega_{\mathbf{p},l} = \emptyset, \end{cases}$$

- $\Omega_{\mathbf{p},l} = \{i; c_i(\mathbf{p} + \mathbf{w}(\mathbf{p})) = l\}$  where  $l=1, \dots, L$  indicates the index set of the voting candidates whose label is  $l$  after being warped to pixel  $\mathbf{p}$ .
- $\tau$  is set to be the value of the maximum difference of SIFT feature:  $\tau = \max_{s_1, s_2, \mathbf{p}} \|s_1(\mathbf{p}) - s_2(\mathbf{p})\|$

# Label Transfer

- Prior term :

$$\lambda(c(\mathbf{p}) = l) = -\log \text{hist}_l(\mathbf{p})$$

- The prior probability that the object category  $l$  appears at pixel  $\mathbf{p}$ .
  - obtained by counting the occurrence of each object category at each location in the training set
  - Location prior

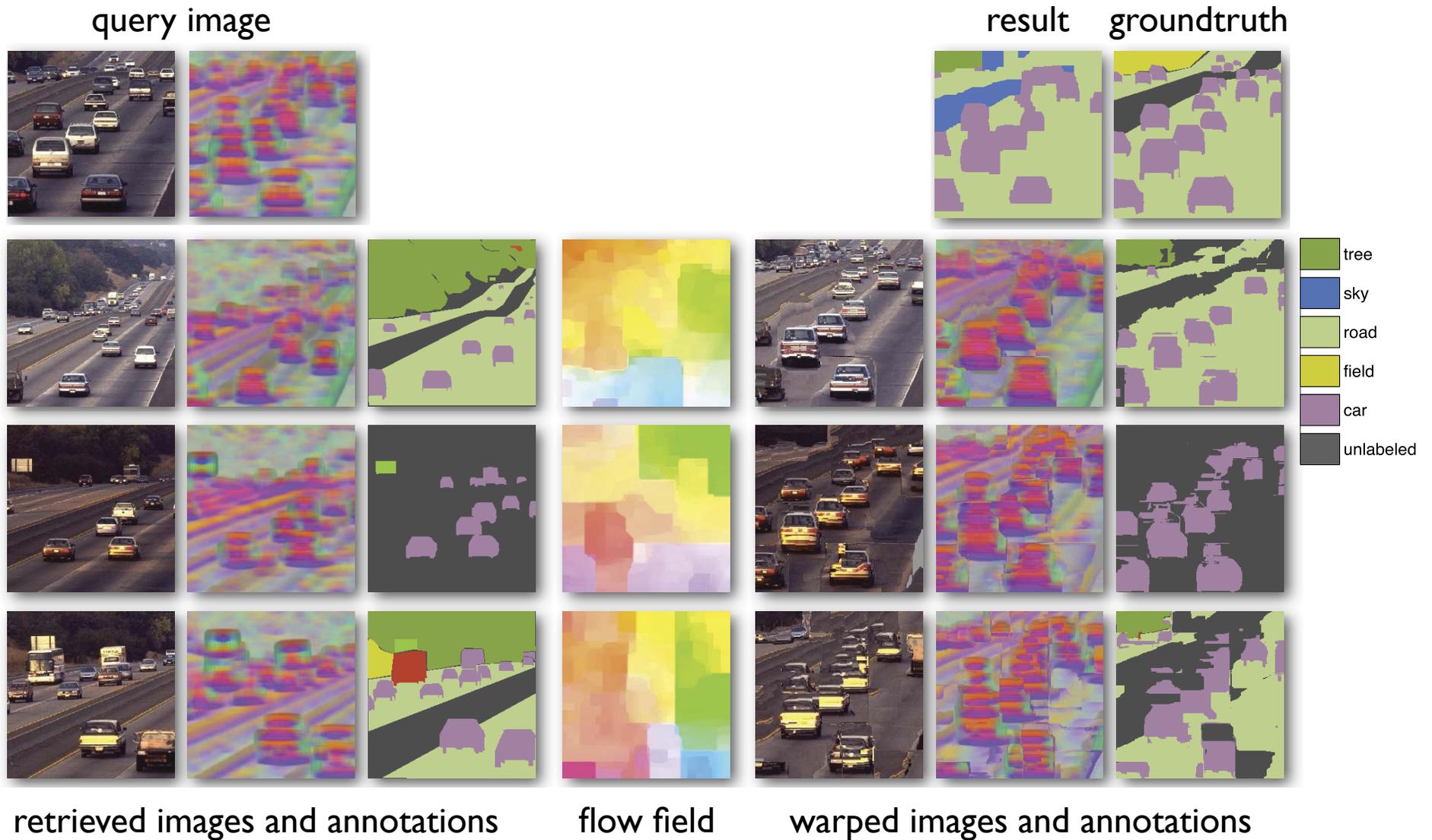
# Label Transfer

- Spatial smoothness term:

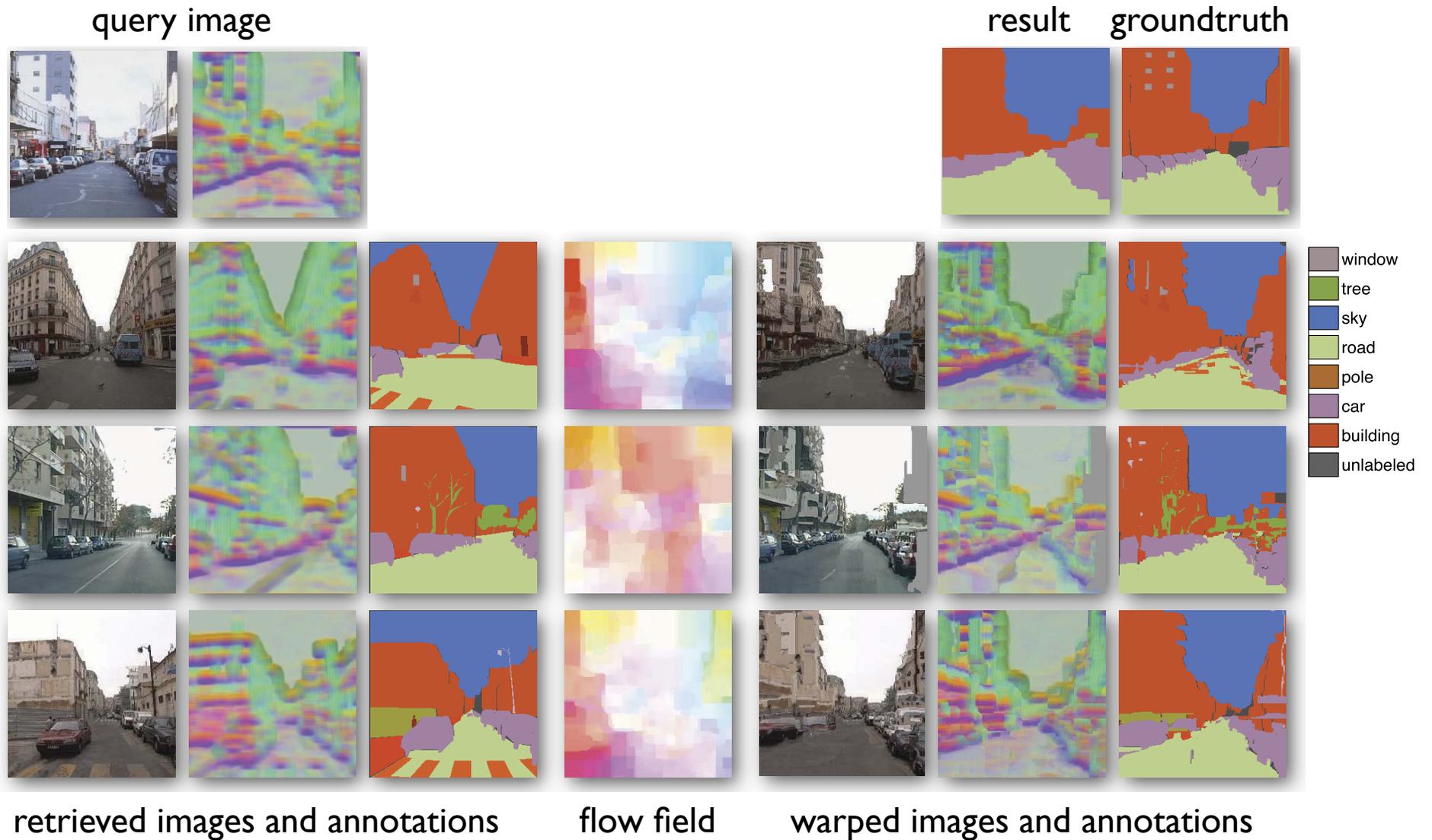
$$\phi(c(\mathbf{p}), c(\mathbf{q})) = \delta[c(\mathbf{p}) \neq c(\mathbf{q})] \left( \frac{\xi + e^{-\gamma \|I(\mathbf{p}) - I(\mathbf{q})\|^2}}{\xi + 1} \right)$$

- The neighboring pixels into having the same label with the probability depending on the image edges:
  - Stronger the contrast, the more likely it is that the neighboring pixels may have different labels.

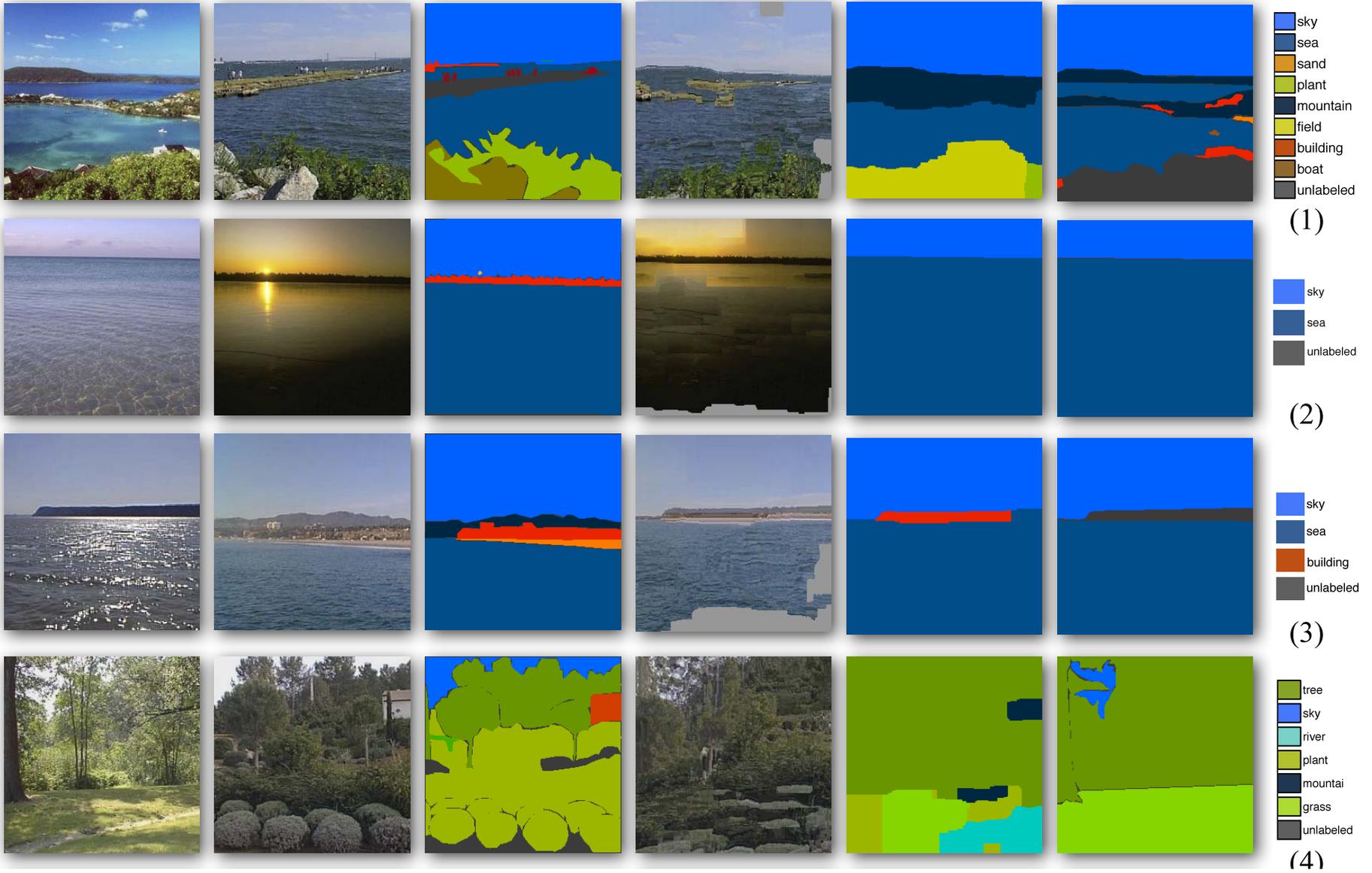
# Parsing Results



# Parsing Results

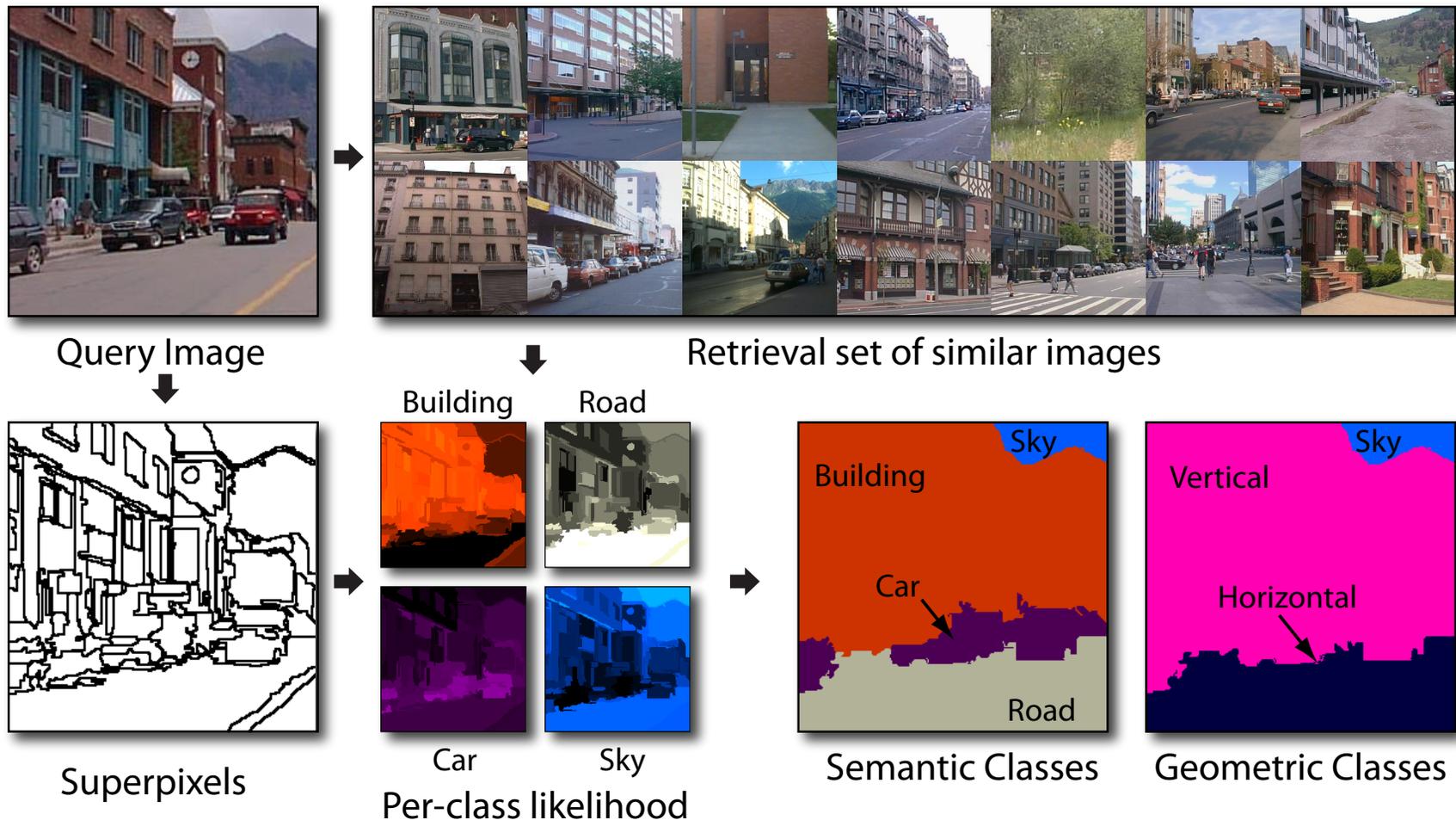


# Parsing Results



# Superparsing: Scalable Nonparametric Image Parsing with Superpixels

(Tighe & Lazebnik, ECCV 2010)



# Superparsing: Scalable Nonparametric Image Parsing with Superpixels

- Framework consists of three main modules:
  1. Scene retrieval: finding nearest neighbor images (k-NN approach)
  2. Superpixel segmentation: segmenting query image into superpixels
  3. Superpixel retrieval: finding nearest neighbors superpixels
  4. Labeling: using a MRF model to label input image

# Scene Retrieval

- Find a relatively small retrieval set of training images that will serve as the source of candidate superpixel-level matches.

(a) Global features for retrieval set computation (Section 2.1)		
Type	Name	Dimension
Global	Spatial pyramid (3 levels, SIFT dictionary of size 200)	4200
	Gist (3-channel RGB, 3 scales with 8, 8, & 4 orientations)	960
	Color histogram (3-channel RGB, 8 bins per channel)	24

# Superpixel Segmentation

- superpixels using the fast graph- based segmentation algorithm of Felzenszwalb and Huttenlocher

(b) Superpixel features (Section 2.2)		
Shape	Mask of superpixel shape over its bounding box ( $8 \times 8$ )	64
	Bounding box width/height relative to image width/height	2
	Superpixel area relative to the area of the image	1
Location	Mask of superpixel shape over the image	64
	Top height of bounding box relative to image height	1
Texture/SIFT	Texton histogram, dilated by 10 pix texton histogram	$100 \times 2$
	Quantized SIFT histogram, dilated by 10 pix quantized SIFT histogram	$100 \times 2$
	Left/right/top/bottom boundary quantized SIFT histogram	$100 \times 4$
Color	RGB color mean and std. dev.	$3 \times 2$
	Color histogram (RGB, 11 bins per channel), dilated by 10 pix color histogram	$33 \times 2$
Appearance	Color thumbnail ( $8 \times 8$ )	192
	Masked color thumbnail	192
	Grayscale gist over superpixel bounding box	320

# Labeling

- A probabilistic MRF model is built to integrate
  - multiple category labels,
  - contextual relations between category labels

$$J(\mathbf{c}) = \sum_{s_i \in SP} E_{\text{data}}(s_i, c_i) + \lambda \sum_{(s_i, s_j) \in A} E_{\text{smooth}}(c_i, c_j)$$

- Superpixel-based labeling as opposed to pixel-level one.

# Data Term

- Likelihood term:

$$E_{\text{data}}(s_i, c_i) = -w_i \sigma(L(s_i, c_i))$$

with  $\sigma(t) = \exp(\gamma t) / (1 + \exp(\gamma t))$  being a sigmoid function.

- constructed from the superpixel matches:

$$\begin{aligned} L(s_i, c) &= \log \frac{P(s_i | c)}{P(s_i | \bar{c})} = \log \prod_k \frac{P(f_i^k | c)}{P(f_i^k | \bar{c})} & \frac{P(f_i^k | c)}{P(f_i^k | \bar{c})} &= \frac{(n(c, \mathcal{N}_i^k) + \epsilon) / n(c, \mathcal{D})}{(n(\bar{c}, \mathcal{N}_i^k) + \epsilon) / n(\bar{c}, \mathcal{D})} \\ &= \sum_k \log \frac{P(f_i^k | c)}{P(f_i^k | \bar{c})}, & &= \frac{n(c, \mathcal{N}_i^k) + \epsilon}{n(\bar{c}, \mathcal{N}_i^k) + \epsilon} \times \frac{n(\bar{c}, \mathcal{D})}{n(c, \mathcal{D})}, \end{aligned}$$

Nonparametric density estimation

# Smoothness Term

- Contextual inference:

$$E_{\text{smooth}}(c_i, c_j) = -\log[(P(c_i|c_j) + P(c_j|c_i))/2] \\ \times \delta[c_i \neq c_j],$$

- enforce contextual constraints.
- defined based on probabilities of label co-occurrence.

# Parsing Results

(a) Query

(b) Ground Truth Labels

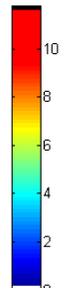
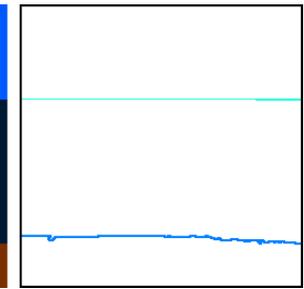
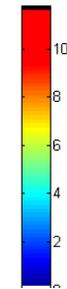
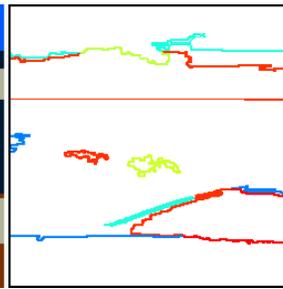
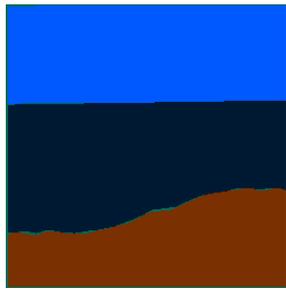
(c) Initial Labeling

(d) Initial Edge Penalties

(e) MRF Labeling

(f) Final Edge Penalties

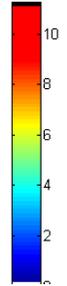
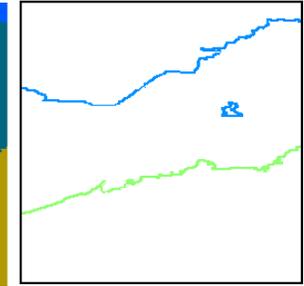
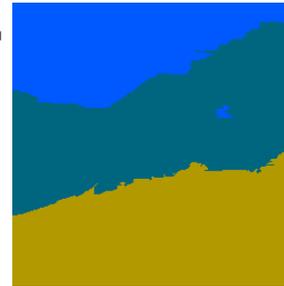
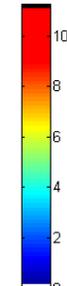
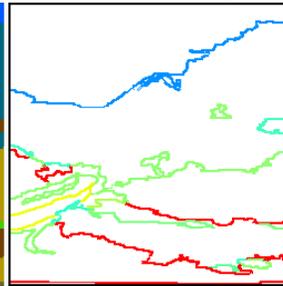
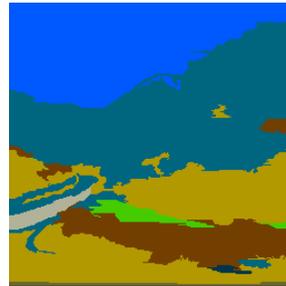
■ Road ■ Sand ■ Sea ■ Sky ■ Tree



72.1

90.2

■ Desert ■ Field ■ Grass ■ Mountain ■ River ■ Road ■ Sky



88.4

93.7

# Including Geometric Priors

- Simultaneous labeling of regions into two types of classes: semantic and geometric

$$H(\mathbf{c}, \mathbf{g}) = J(\mathbf{c}) + J(\mathbf{g}) + \mu \sum_{s_i \in SP} \varphi(c_i, g_i)$$

- Three geometric labels:
  - sky, horizontal, and vertical
- enforces coherence between the geometric and semantic labels.

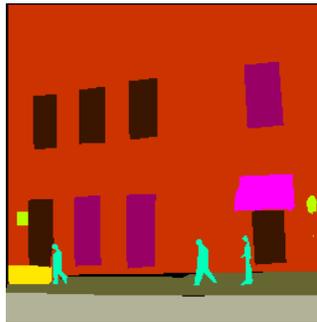
$\varphi(c_i, g_i)$ : a coherence term between the semantic and geometric labels of the same superpixel.

# Parsing Results

(a) Query



(b) Ground Truth Labels

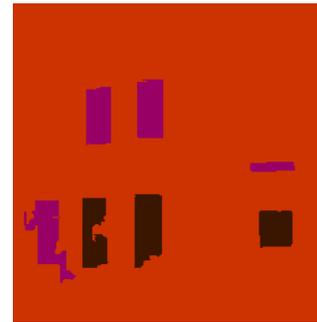


(c) Initial Labeling



53.2

(d) Semantic MRF



67.9

(e) Joint Semantic and Geometric



68.7

- |  |   |  |
|--|---|--|
|  Awning     |  Road        |  Sky    |
|  Balcony  |  Sidewalk  |  Vert |
|  Building |  Sign      |  Horz |
|  Door     |  Staircase |  |
|  Person   |  Window    |  |



97.8



97.6



97.6

# Parsing Results



# Parsing Results

Query

Ground Truth Labels

Initial Labeling

Final Labeling

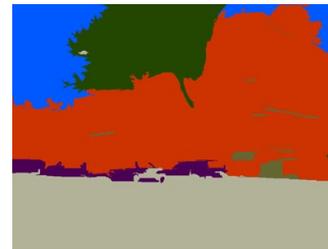
Geometric Labeling



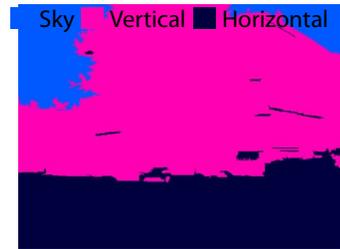
(e)



72.3



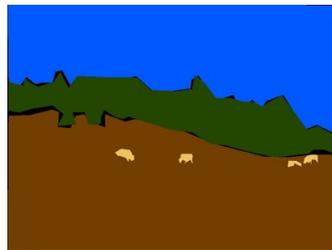
79.1



90.3



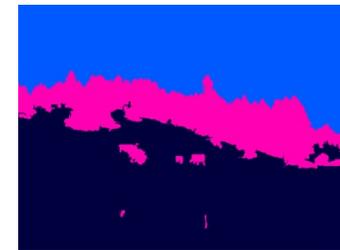
(f)



95.1



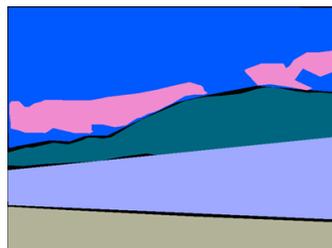
96.4



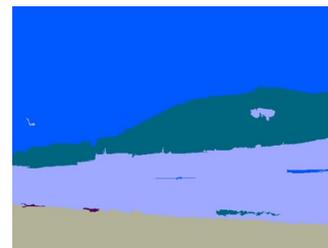
96.8



(g)



67.7



89.5



98.5

