# BIL 717 Image Processing

## **Semantic Segmentation**

Erkut Erdem Hacettepe University Computer Vision Lab (HUCVL)

# Review - Solving MRFs with graph cuts

#### Main idea:

- Construct a graph such that every st-cut corresponds to a joint assignment to the variables y
- The cost of the cut should be equal to the energy of the assignment, E(y; data)\*.
- The minimum-cut then corresponds to the minimum energy assignment, y<sup>\*</sup> = argmin<sub>y</sub> E(y; data).





\* Requires non-negative energies

S. Gould







# Review - Why Higher-order Functions?

In general  $\theta(x_1, x_2, x_3) \neq \theta(x_1, x_2) + \theta(x_1, x_3) + \theta(x_2, x_3)$ 

#### 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]



MRF with

 $\mathsf{E}(\mathsf{x}) = \sum \Theta_{ii} (\mathsf{x}_i, \mathsf{x}_i)$ 

i.i € №

global variables

Order 2

Higher-order MRF

 $+\theta(x_1,\ldots,x_n)$ 

C. Rother

Order n

 $E(\mathbf{x}) = \sum \Theta_{ii} (\mathbf{x}_i, \mathbf{x}_i)$ 

II C N

# **Semantic Segmentation**

• Joint recognition & segmentation

higher(8)-connected;

 $\mathsf{E}(\mathsf{x}) = \sum \Theta_{ij} \left( \mathsf{x}_{i}, \mathsf{x}_{j} \right)$ 

Order 2

i.i € N₀

pairwise MRF

4-connected

pairwise MRF

iien.

Order 2

 $\mathsf{E}(\mathsf{x}) = \sum \Theta_{ij} \left( \mathsf{x}_{i}, \mathsf{x}_{j} \right)$ 

- 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

M. P. Kumar

# Early Studies of Semantic Segmentation



# Early Studies of Semantic Segmentation



R. Fergus

# Early Studies of Semantic Segmentation

Using Normalized Cuts, Shi & Malik, 1997





# Jigsaw approach: Borenstein and Ullman, 2002



# **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





C. Rother

skv

arass

[TextonBoost; Shotton et al, '06]

[TextonBoost; Shotton et al, '06]





- Framework consists of three main modules:
  - 1. Scene retrieval: finding nearest neighbors (k-NN approach)
  - 2. Dense scene alignment: dense scene matching (SIFT Flow)



### 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



# Label Transfer

- A set of voting candidates {s<sub>i</sub>, c<sub>i</sub>, w<sub>i</sub>}<sub>i=1:M</sub> is obtained from the retrieved images with s<sub>i</sub>, c<sub>i</sub>, and w<sub>i</sub> 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

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

### 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

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

$$\sum_{\mathbf{p}} \eta(|u(\mathbf{p})| + |v(\mathbf{p})|) +$$

small displacement term

$$\begin{split} & \sum_{(\mathbf{p},\mathbf{q})\in\varepsilon}\min(\lambda|u(\mathbf{p})-u(\mathbf{q})|,d) + \\ & \min(\lambda|v(\mathbf{p})-v(\mathbf{q})|,d), \end{split}$$

smoothness term

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

# Label Transfer

• Likelihood term:

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

- Ω<sub>p,l</sub> = {i; c<sub>i</sub>(p + w(p)) = l} where l=1,...,L indicates the index set of the voting candidates whose label is l after being warped to pixel 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 :

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

- The prior probability that the object category *I* appears at pixel **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**



