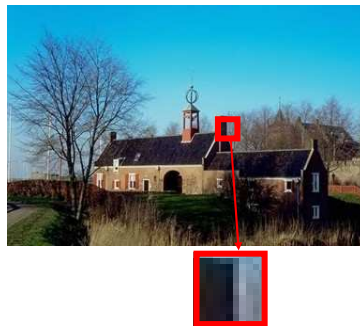


**Martial Hebert  
Sanjiv Kumar**

**The Robotics Institute  
School of Computer Science  
Carnegie Mellon University**

1

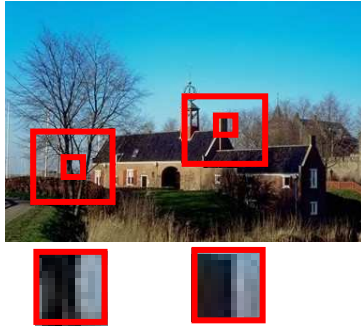
## Motivation



**Map to appropriate feature space and classify**

2

## Motivation



**Need context from larger neighborhoods !**

3



**Context from whole image !**

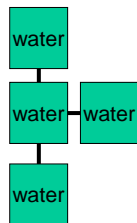
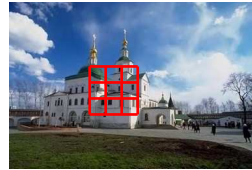
4



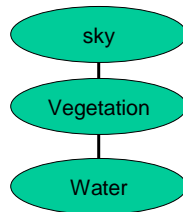
Context from Other Objects !

5

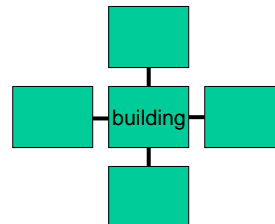
## Interactions



Pixel-Pixel



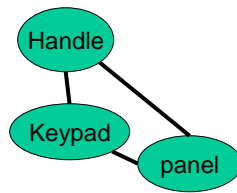
Region-Region



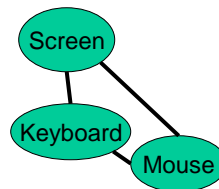
Patch-Patch

6

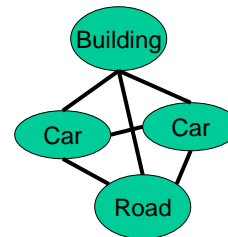
## Interactions



Part-Part



Object-Object



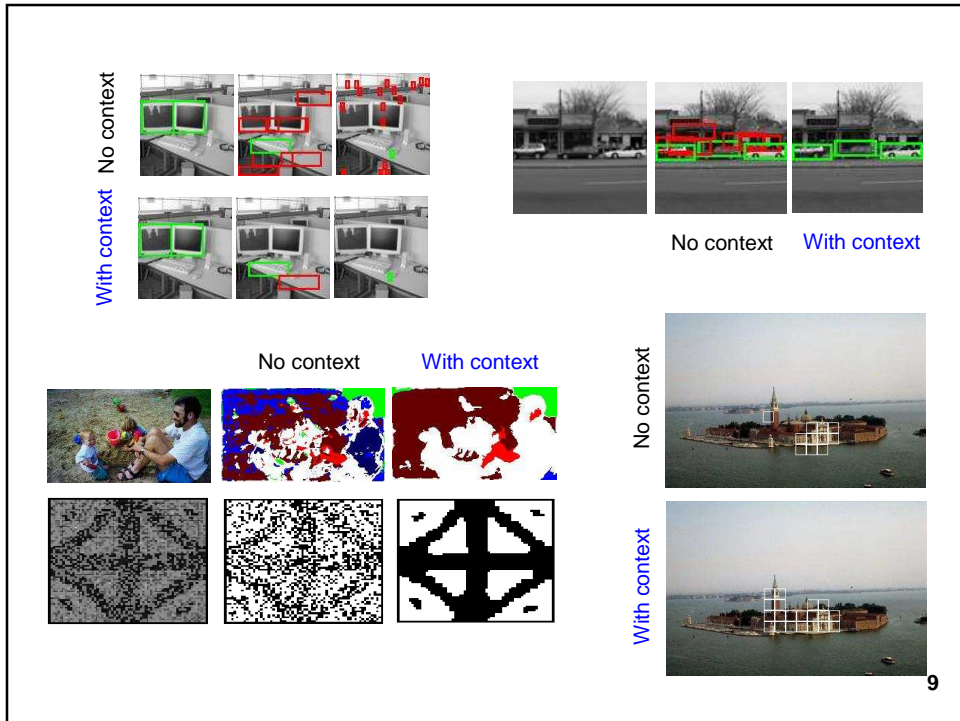
Region-Object

7

## Modeling Interactions

- Framework to learn **all** relevant contextual interactions in a **single model** automatically from training data
- Probabilistic models
  - Principled way to deal with ambiguities
- Graphical models
  - Powerful framework for ensuring global consistency using relatively local constraints

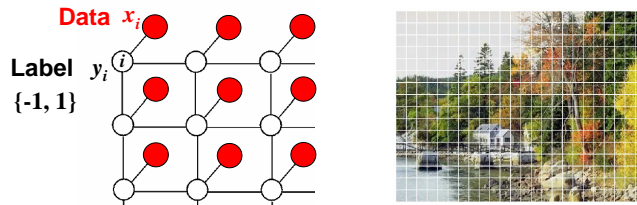
8



## Outline

- **Background**
  - Markov Random Fields (MRFs)
- **Conditional Random Fields (CRFs)**
- **Multiclass**
- **Hierarchical Interactions**

## Markov Random Field (MRF)



Generative Framework

$$P(y|x) \propto P(y, x) = \underbrace{P(x|y)}_{\text{Observation Model}} \underbrace{P(y)}_{\text{Prior Model (MRF)}}$$

Observation Model Prior Model (MRF)

$$P(y|x) = \frac{1}{Z} \exp \left( \sum_{i \in S} \log P(x_i | y_i) + \sum_{i \in S} \sum_{j \in N_i} \beta y_i y_j \right)$$

[Geman & Geman, '84]

11

## Generative vs. Discriminative

We want :  $P(y|x)$

- Generative Framework
  - Models  $P(y, x)$  to get  $P(y|x)$
  - Implicit modeling of observations
- Discriminative Framework
  - Models Conditional distribution  $P(y|x)$  directly

12

## Conditional Random Field (CRF)

- $P(y|x)$  is directly modeled as Gibbs Field
- Graphs on images with loops
- Clique potentials using local discriminative models

$$P(y|x) = \frac{1}{Z} \exp \left( \underbrace{\sum_{i \in S} A_i(y_i, x)}_{\text{Association Potential}} + \underbrace{\sum_{i \in S} \sum_{j \in N_i} I_{ij}(y_i, y_j, x)}_{\text{Interaction Potential}} \right)$$

[Lafferty et al., ICML'01][Kumar, ICCV '03]

13

## Comparison with MRF framework

$$\begin{array}{l} \text{MRF} \quad P(y|x) = \frac{1}{Z} \exp \left( \sum_{i \in S} \log P(x_i | y_i) + \sum_{i \in S} \sum_{j \in N_i} \beta y_i y_j \right) \\ \text{CRF} \quad P(y|x) = \frac{1}{Z} \exp \left( \sum_{i \in S} A(y_i, x) + \sum_{i \in S} \sum_{j \in N_i} I(y_i, y_j, x) \right) \end{array}$$

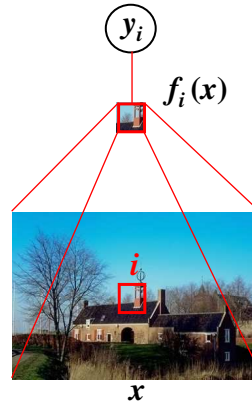
↑ ↑  
Data from multiple sites Data-dependent label interactions

14

## Example Association Potential $A(y_i, x)$

discriminative classifier

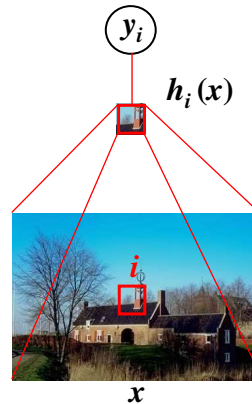
$$\begin{aligned} A(y_i, x) &= \log P(y_i | f_i(x)) \\ &= \log \sigma(y_i w^T f_i(x)) \end{aligned}$$



## Example Association Potential $A(x_i, y)$

discriminative classifier

$$\begin{aligned} A(y_i, x) &= \log P(y_i | h_i(x)) \\ &= \log \sigma(y_i w^T h_i(x)) \end{aligned}$$



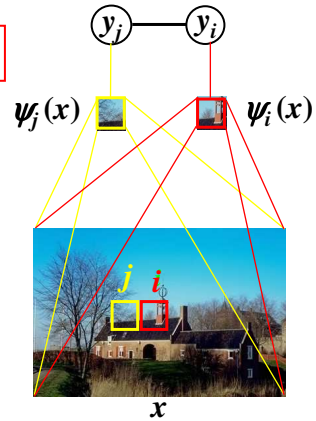
Mapping to induce a nonlinear decision boundary:

$$h_i(x) = [1, \phi_1(f_i(x)), \dots, \phi_R(f_i(x))]^T$$

## Example Interaction Potential $I(y_i, y_j, x)$

pairwise discriminative classifier

$$\begin{aligned} I(y_i, y_j, x) &= \log P(y_i, y_j | \psi_i, \psi_j) \\ &= y_i y_j \underbrace{v^T \mu_{ij}(x)}_{\beta} \end{aligned}$$



[Kumar & Hebert, NIPS '04]

17

## Using the Model

### *Parameter Learning*

$M$  training images  $\rightarrow$  parameters  $\theta = \{w, v\}$

### *Inference*

Input image  $\rightarrow$  Labels at every site  $y_i$

18

## Parameter Learning

- Maximum Likelihood - Given  $M$  training images,

$$l(\theta) = \sum_{m=1}^M \left\{ \sum_{i \in S} \log \sigma(y_i \mathbf{w}^T \mathbf{h}_i(\mathbf{x})) + \sum_{i \in S} \sum_{j \in N_i} y_i y_j v^T \mu_{ij}(\mathbf{x}) - \log \mathbf{Z} \right\}$$

$$\frac{\partial l(\theta)}{\partial \mathbf{w}} = \frac{1}{2} \sum_m \sum_i (y_i - \langle y_i \rangle_{P(y|x)}) \mathbf{h}_i(\mathbf{x})$$

$$\frac{\partial l(\theta)}{\partial \mathbf{v}} = \sum_m \sum_i \sum_j (y_i y_j - \langle y_i y_j \rangle_{P(y|x)}) \mu_{ij}(\mathbf{x})$$

- Maximum pseudo-likelihood
- Contrastive Divergence
- Marginal Approximations
- Saddle Point Approximation

[Hinton '86]

19

## Parameter Learning

- Maximum pseudo-likelihood (PL)

$$\hat{\theta} = \arg \max_{\theta} \prod_{m=1}^M \prod_{i \in S} P(y_i^m | y_{N_i}^m, \mathbf{x}^m, \theta)$$

- Overestimates interaction parameters [Kumar & Hebert '03]

- Contrastive Divergence (CD) [Hinton '02]

$$\langle y_i \rangle \approx (y_i)_{P_T^1} \quad \text{and} \quad \langle y_i y_j \rangle \approx (y_i y_j)_{P_T^1}$$

- Large bias [Williams et al. '02]

20

## Marginals Based Approximation

- Pseudo Marginal Approximation (PMA)

$$\langle y_i \rangle = \sum_{y_i} y_i P(y_i | x) \approx \sum_{y_i} y_i b_i(y_i)$$

Pseudo marginals (Belief Propagation)

$$\langle y_i y_j \rangle = \sum_{y_i} \sum_{y_j} y_i y_j P(y_i, y_j | x) \approx \sum_{y_i} \sum_{y_j} y_i y_j b_{ij}(y_i, y_j)$$

- Maximum Marginal Approximation (MMA)

$$\tilde{y}_i = \arg \max_{y_i} b_i(y_i)$$

$$\langle y_i \rangle \approx \tilde{y}_i \text{ and } \langle y_i y_j \rangle \approx \tilde{y}_i \tilde{y}_j$$

21

## Saddle Point Approximation (SPA)

- Maximum Likelihood - Given  $M$  training images,

$$\frac{\partial \mathcal{L}(\theta)}{\partial \mathbf{w}} = \frac{1}{2} \sum_m \sum_i (y_i - \langle y_i \rangle_{P(y|x)}) h_i(\mathbf{x})$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \mathbf{v}} = \sum_m \sum_i \sum_j (y_i y_j - \langle y_i y_j \rangle_{P(y|x)}) \mu_{ij}(\mathbf{x})$$

- Saddle Point Approximation

$$\hat{y} = \arg \max_y P(y | \mathbf{x}, \theta)$$

$$\langle y_i \rangle \approx \hat{y}_i, \quad \langle y_i y_j \rangle \approx \hat{y}_i \hat{y}_j$$

22

## Saddle Point Approximation (SPA)

$$\hat{y} = \arg \max_y P(y|x, \theta)$$

$$\frac{\partial l(\theta)}{\partial w} \approx \frac{1}{2} \sum_m \sum_i (y_i - \hat{y}_i) h_i(x)$$

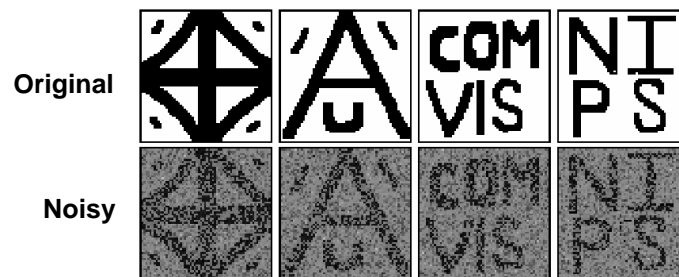
$$\frac{\partial l(\theta)}{\partial v} \approx \sum_m \sum_i \sum_j (y_i y_j - \hat{y}_i \hat{y}_j) \mu_{ij}(x)$$

Perceptron Learning Rules

[Freund et al. ML'99] [Collins, EMNLP'02]  
[LeCun et al. IEEE Proc.'98]

23

## Comparative Experiments



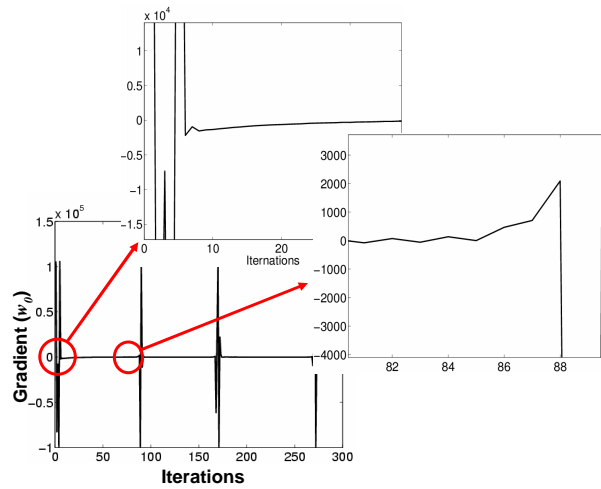
- 4 base images (64 x 64)
- Corrupted by two noise models
  - Gaussian (single Gaussian)
  - Bimodal (mixture of two Gaussians)
- Training – 10 images, Testing – 200 images

24

## Approximate Gradient Ascent



QuickTime™ and a decompressor are needed to see this picture.



[Kumar & Jonas, EMMCVPR '05]

25

## Inference

Given input image  $x$  and  $P(y|x)$ , get the optimal labels  $y$

- **Exact Maximum A Posteriori (MAP)**
  - Mincut-maxflow [ Greig et. al '89 ]
- **Approximate Maximum Marginals**
  - Loopy Belief Propagation (BP) [ Frey & Mckay '98 ]

26

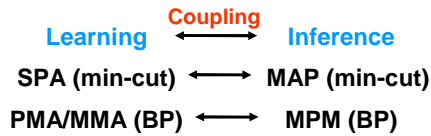
## Learning Experiments / Bimodal Noise

Inference

Parameter Learning

	MAP (min-cut)	MPM (BP)	Local MAP (ICM)
SPA (min-cut)	<b>5.82</b>	19.19	14.88
PMA (BP)	6.45	<b>5.48</b>	17.39
MMA (BP)	26.53	<b>5.70</b>	16.00
CD	8.98	6.29	8.91
PL	17.69	7.30	22.21

Pixelwise error (%) on 200 test images



27

[Kumar & Jonas, EMCCVPR '05][Zhu et al., PAMI'02]

## Example: Man-Made Structure Detection

Training Set – 130 images (256x384 pixels)



Scale variations

Illumination variations

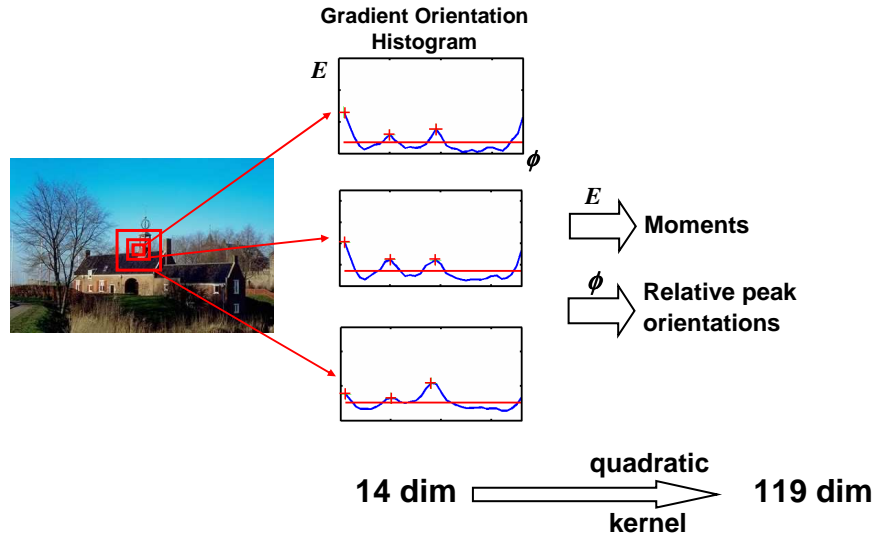
Pose variations

Non-linear structures

Negative samples

28

## Design of Feature Vector ( $h_i$ )

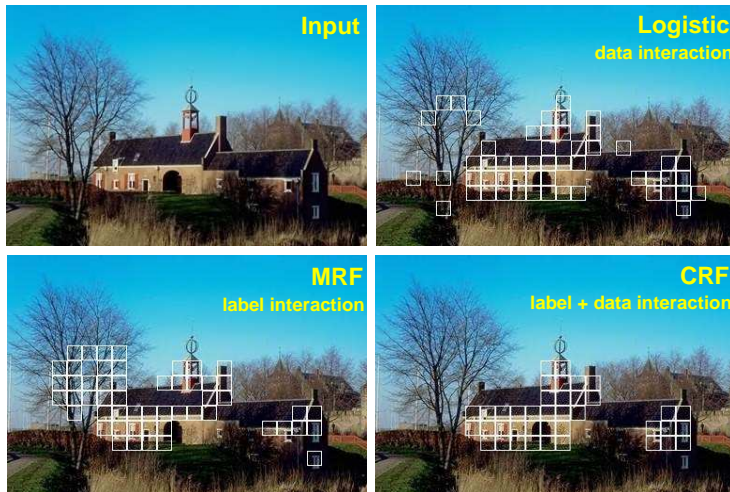


[Kumar & Hebert CVPR'03]

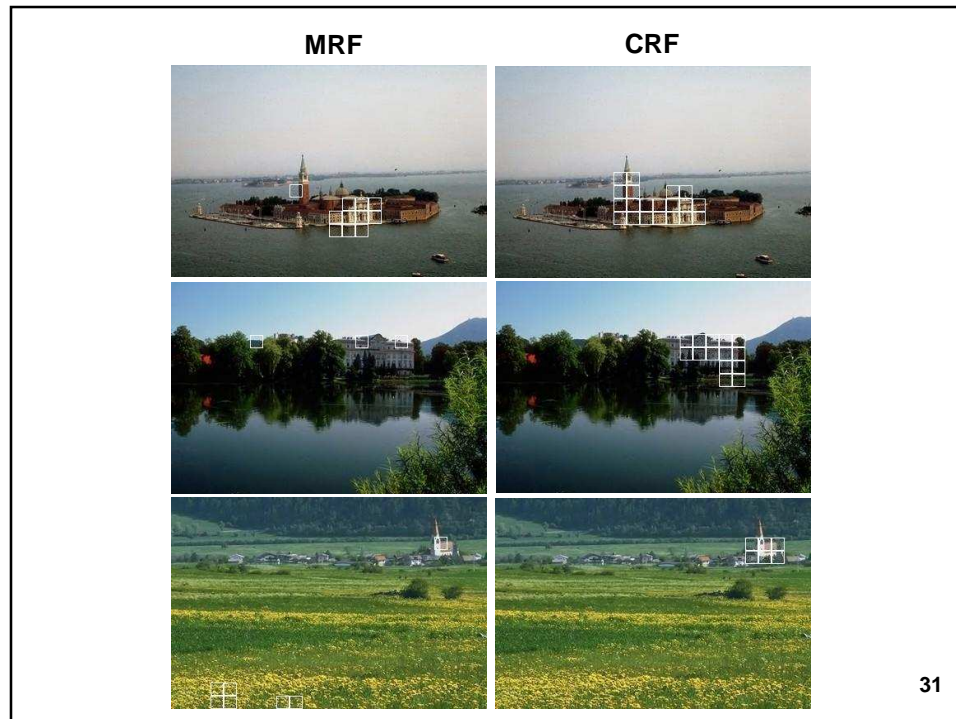
[Lafferty et al., ICML'04]

29

## Performance Evaluation



30



## So Far...

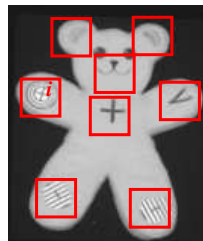
- **Background**
  - Markov Random Fields (MRFs)
- **Conditional Random Fields (CRFs)**
- **Multiclass**
- **Hierarchical Interactions**

## Multiclass Extension

- Labels  $y_i \in \{1, \dots, C\}$
- Association Potential

$$A(y_i, x) = \log P(y_i | x)$$

$$P(y_i = k | x) = \begin{cases} \frac{\exp(\mathbf{w}_k^T \mathbf{h}_i(x))}{1 + \sum_{l=1}^{C-1} \exp(\mathbf{w}_l^T \mathbf{h}_i(x))} & \text{if } k < C \\ \frac{1}{1 + \sum_{l=1}^{C-1} \exp(\mathbf{w}_l^T \mathbf{h}_i(x))} & \text{if } k = C \end{cases}$$



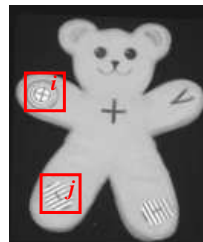
[Kumar & Hebert, ICCV'05][Quattoni et al., NIPS'04]

33

## Multiclass Extension

- Interaction Potential

$$I(y_i, y_j, x) = \sum_{k=1}^C \sum_{l=1}^C \mathbf{v}_{kl}^T \mu_{ij}(x) \delta(y_i = k) \delta(y_j = l)$$



- Inference : Loopy Belief Propagation (BP)
  - Alternatives: CCCP [Yuille '02], EP [Minka '01]

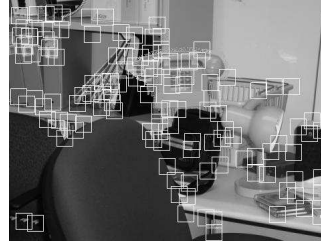
34

## Parts-Based Object Detection

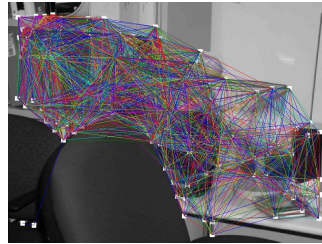
Input image



Patch detection



Graph formation



Detection results



Inference (Loopy BP) time: 1.35 sec

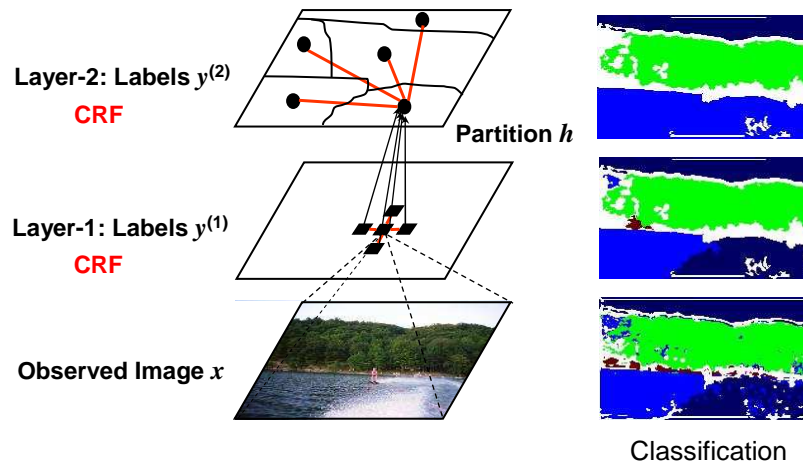
35

## Outline

- **Background**
  - Markov Random Fields (MRFs)
- **Conditional Random Fields (CRFs)**
- **Multiclass**
- **Hierarchical Interactions**

36

## Hierarchical Interactions



37

## Hierarchical Formulation

We want: Labels given input image

Layer-2 labels given mapping and labels 1

Mapping from layer-1 to 2

Layer-1 labels given input image

$$P(y | x) \cong \sum_{h, y^{(1)}} P(y^{(2)} | h, y^{(1)}) P(h | y^{(1)}) P(y^{(1)} | x)$$

CRF

CRF

38

## Design Issues

$$P(y^{(2)} | h, y^{(1)})$$

$$P(y^{(1)} | x)$$

Designed conditional random fields at layer 1 and layer 2

$$\sum_{y^{(1)}}$$

Leads to combinatorial explosion and cannot be computed exactly

Designed approximation (using the fact that  $P(y^{(1)} | x)$  is close to 0 almost everywhere)

$$P(h | y^{(1)})$$

Designed a simple distribution model

39

## (An Attempt at a) Hierarchical Formulation

We want:  
Labels given  
input image

Layer-2 labels  
given mapping  
and labels 1

Mapping from  
layer-1 to 2

Layer-1 labels  
given input  
image

$$P(y | x) \cong \sum_{h, y^{(1)}} P(y^{(2)} | h, y^{(1)}) P(h | y^{(1)}) \underbrace{P(y^{(1)} | x)}_{\delta(y^{(1)} - \tilde{y}_b)}$$

$b$  = Marginals at layer 1  $b_i = P(y_i^{(1)} | x)$

$\tilde{y}_b$  = Max. marginal at layer 1  
(estimated through BP)

40

## (An Attempt at a) Hierarchical Formulation

We want:  
Labels given  
input image
Layer-2 labels  
given mapping  
and labels 1
Mapping from  
layer-1 to 2
Layer-1 labels  
given input  
image

$$P(y | x) \cong \sum_{h, y^{(1)}} P(y^{(2)} | h, y^{(1)}) P(h | y^{(1)}) P(y^{(1)} | x)$$

$\approx \sum_h P(y^{(2)} | h, b) P(h | b)$

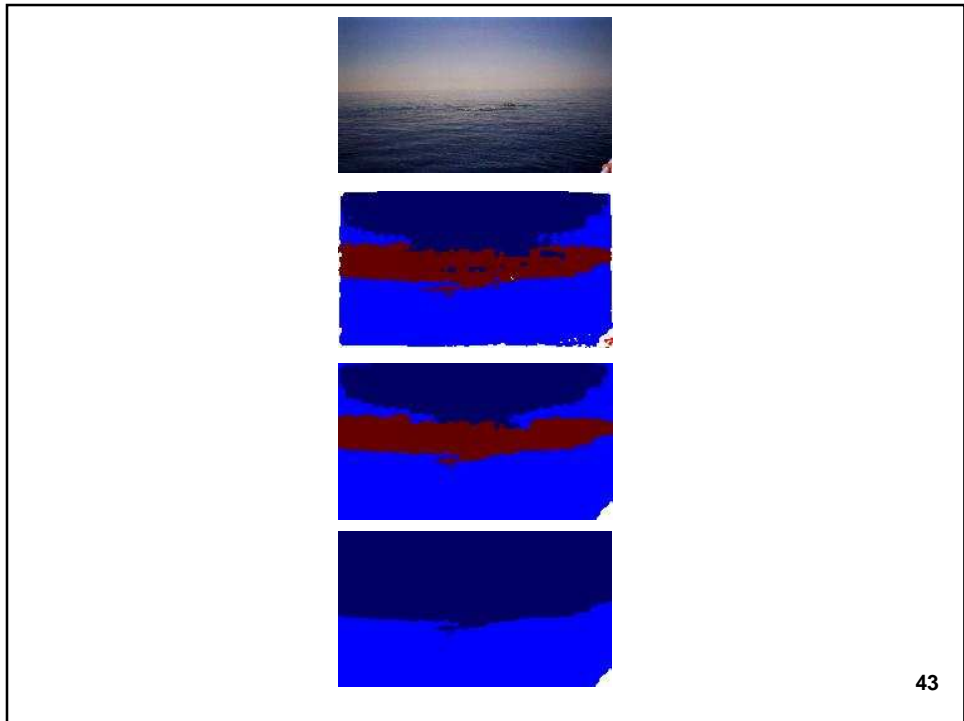
- Sequential Training – Pseudo marginals (BP)
- Inference – Max marginals (BP)

41

## Region Classification

- Data sets - Images of outdoor scenes
  - Beach Dataset: 123 images (124x218 pixels), 6 classes
  - Sowerby Dataset: 104 images (64x96 pixels), 7 classes
- Layer 1 Features
  - HSV/Luv Color, texture moments, oriented DoG filters
- Layer 2 features
  - Relative configuration of regions (*above, beside, enclosed*)
- Training Time: ~ 10 min on 2.8 GHz machine

42



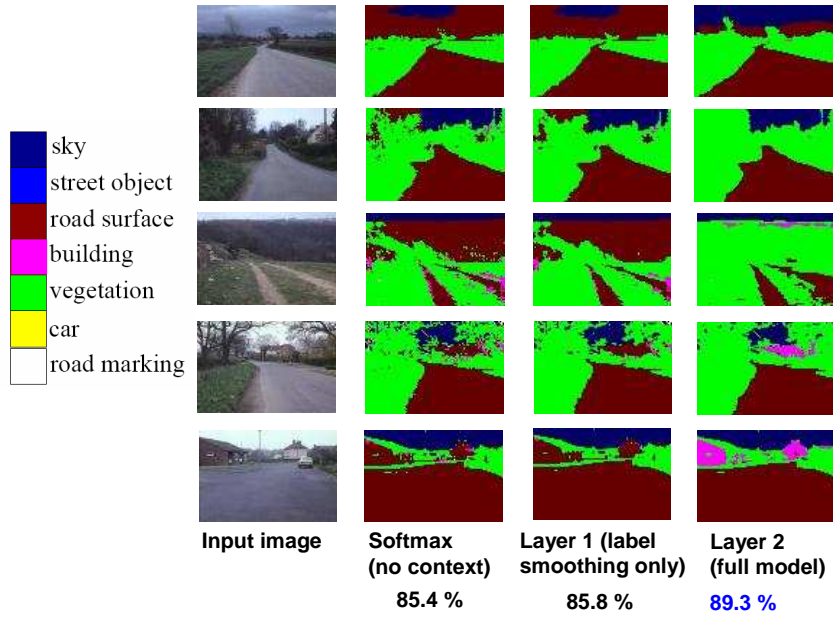
43

### Beach Dataset

	<b>Input image</b>	<b>Softmax (no context)</b> 62.3 %	<b>Layer-1 (label smoothing only)</b> 63.8 %	<b>Layer-2 (full model)</b> <b>74 % (~ 2 Sec)</b>
			<b>Previous best 64 %</b>	<b>44</b>

[Singhal et al., CVPR'04]

## Sowerby Dataset

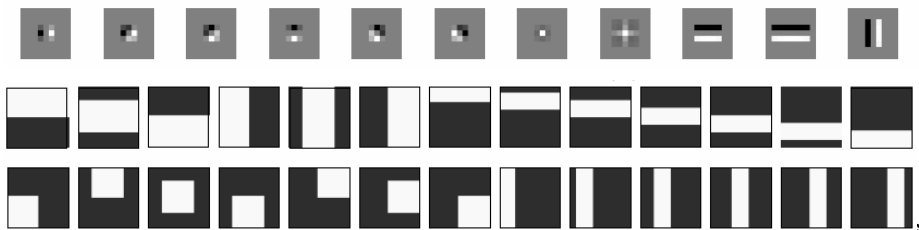


[He et al., CVPR'04][Feng et al., PAMI'02]

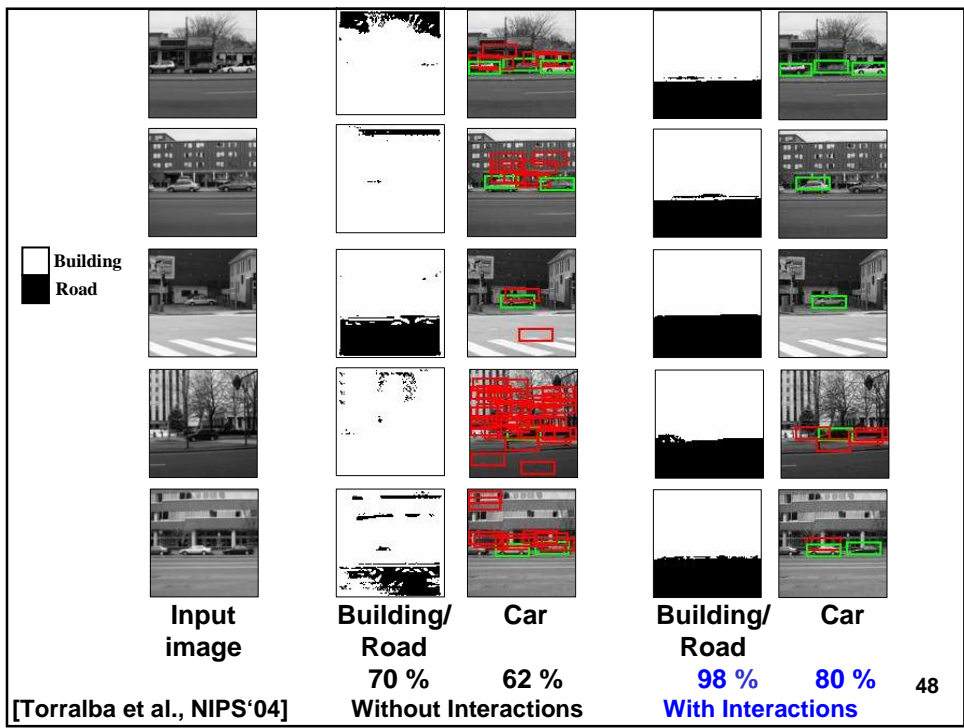
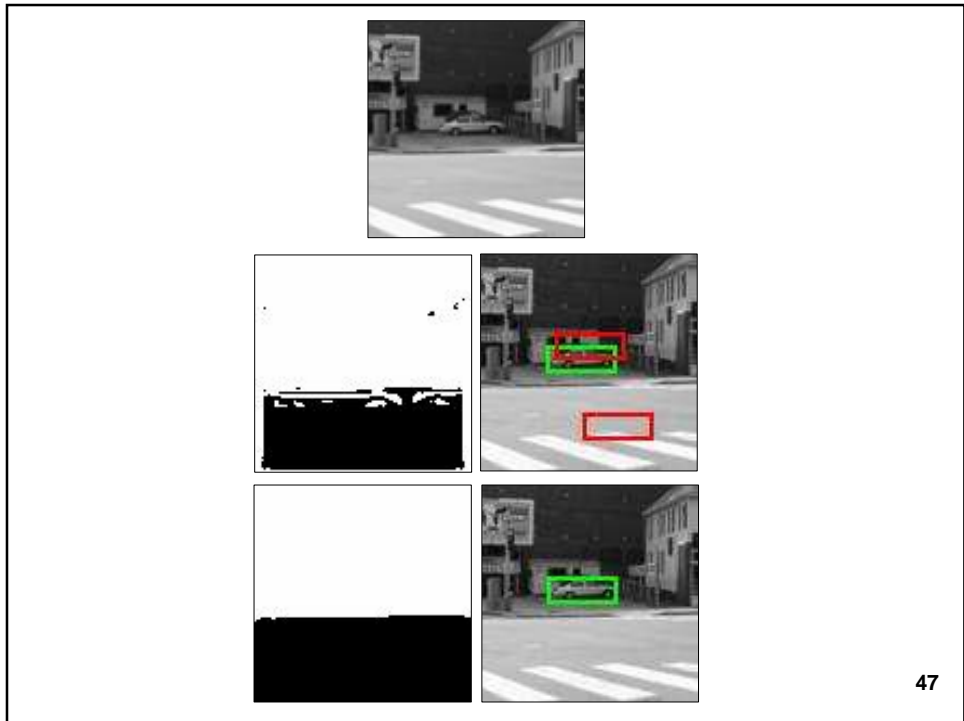
45

## Object-Region Interactions

- **MIT database**  
~100 outdoor images cars, buildings and roads (100x100 pixels)
  - **Region features**
    - Intensity, oriented DoG filters for texture
  - **Object features**
    - Detector score trained with Gentle Boosting
- Dictionary (500 filters and templates)

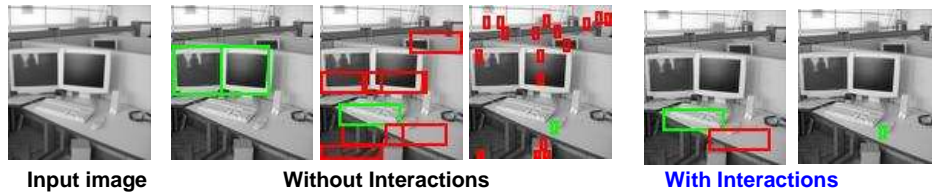


[Torralba et al., NIPS'04]

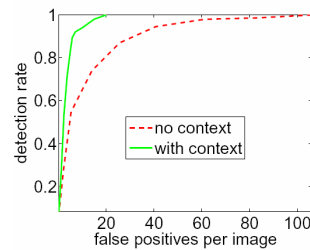


## Object-Object Interactions

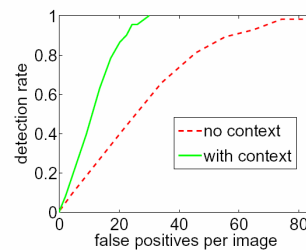
- MIT context database: 164 images (100x100 pixels)
- Very small objects (8x5 pixels) → High false positives
- Initial object detectors trained with Gentle Boosting



Keyboard



Mouse



49

## Summary

- Background
  - Markov Random Fields (MRFs)
- Conditional Random Fields (CRFs)
- Multiclass
- Hierarchical Interactions

50