

# Context and Scene Parsing

Computer Vision

CS 143, Brown

James Hays

Many Slides from  
Svetlana Lazebnik

# Recap: Context and Spatial Layout

- Contextual Reasoning: making a decision based on more than *local* image evidence.
- Numerous sources of context can be exploited to improve scene understanding
- We discussed spatial layout in particular
  - “Geometric Context” method of Hoiem et al.
  - Geometry as a single view *recognition* problem, rather than a multi-view problem.

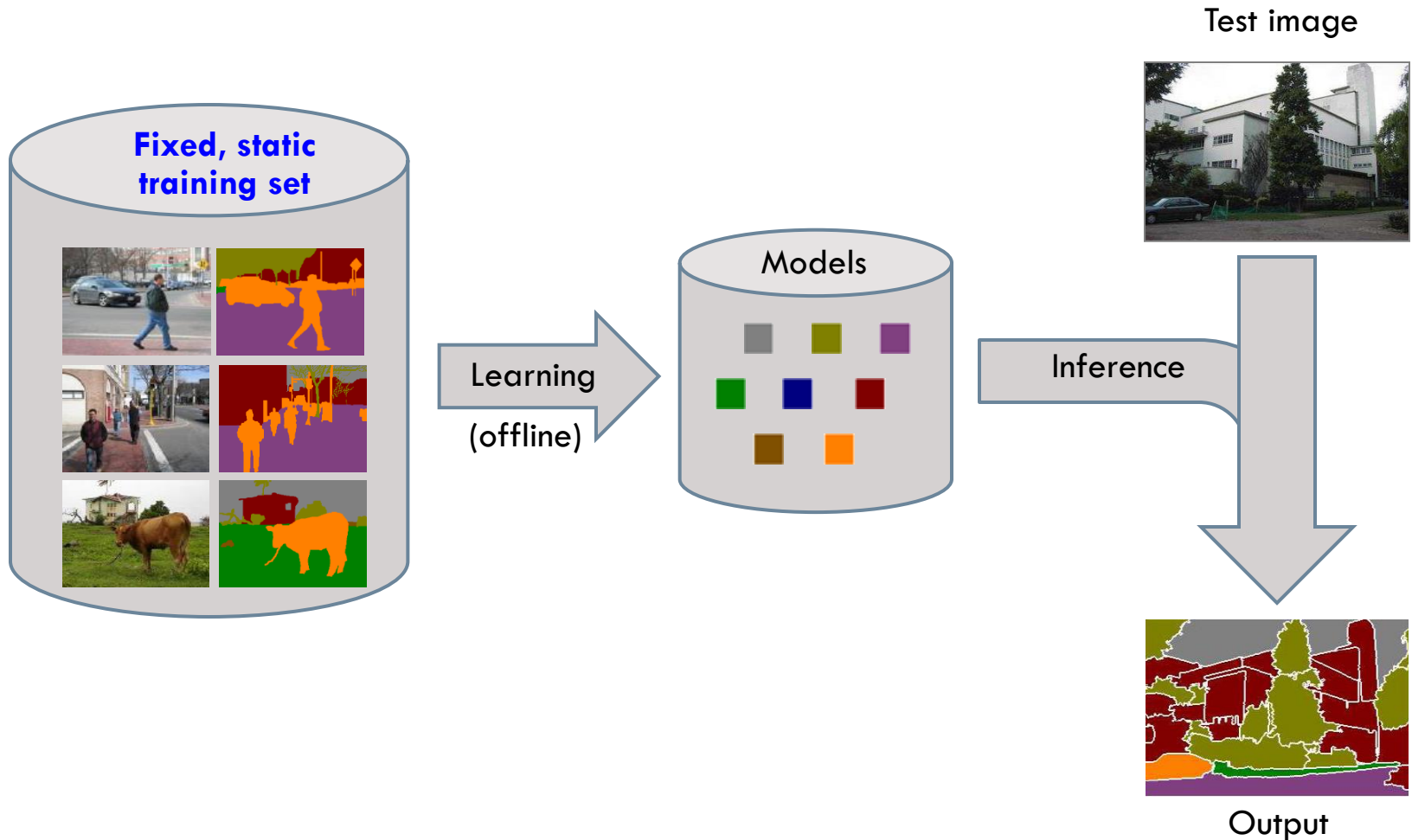
# Today: Scene Parsing

- Label every pixel of an image with a category label (usually with the help of contextual reasoning).
- We'll look at the “non parametric” approach of Tighe and Lazebnik

# Closed-universe recognition

**Fixed, pre-defined set of classes**

■ sky ■ tree ■ road ■ grass ■ water ■ bldg ■ mntn ■ fg obj.



# Closed-universe datasets



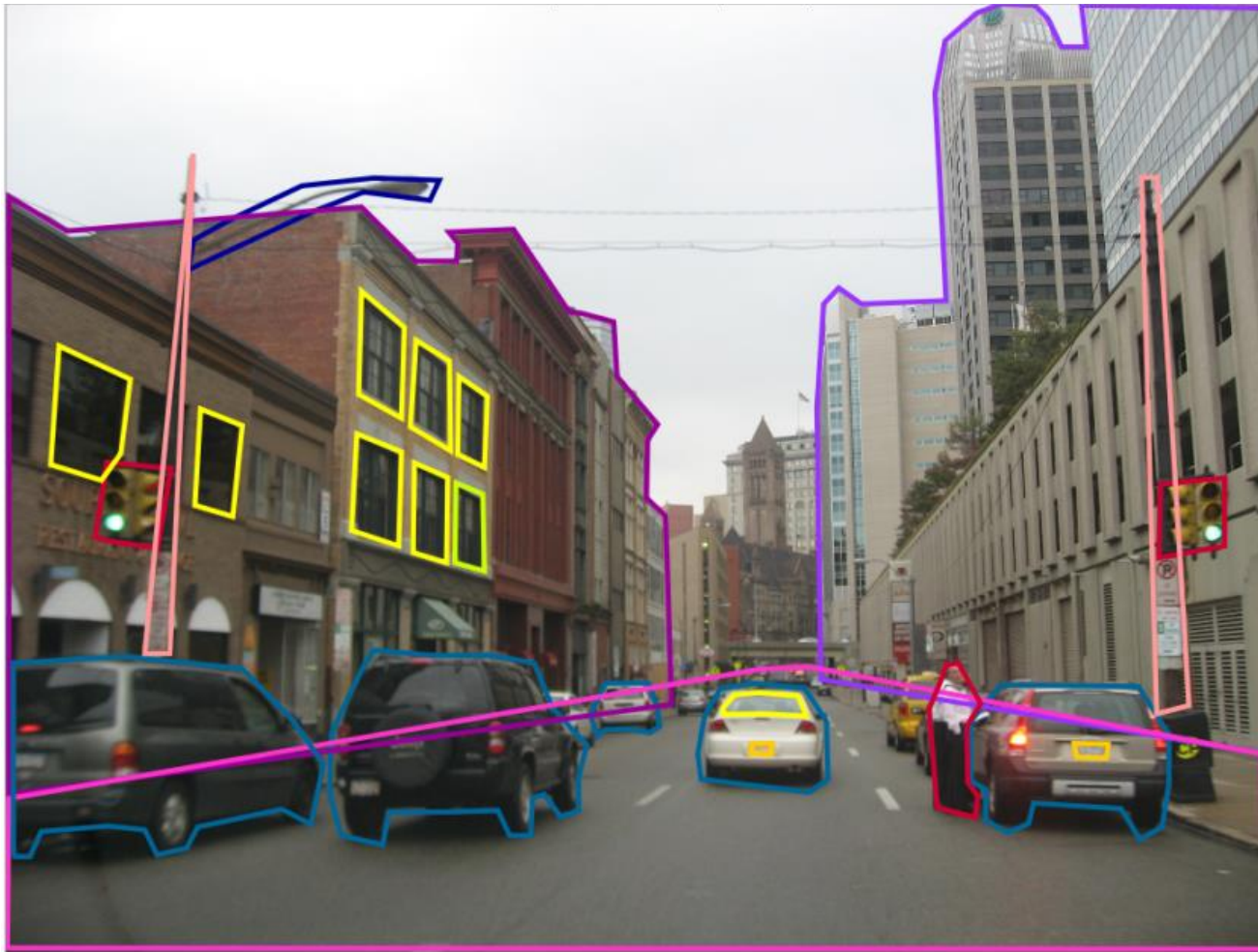
- Small amount of data
- Static datasets
- Limited variation
- Full annotation

# Open-universe datasets



- Large amount of data
- Evolving datasets
- Wide variation
- Incomplete annotation

# Open-universe recognition



There are **754152** labelled objects

## Polygons in this image

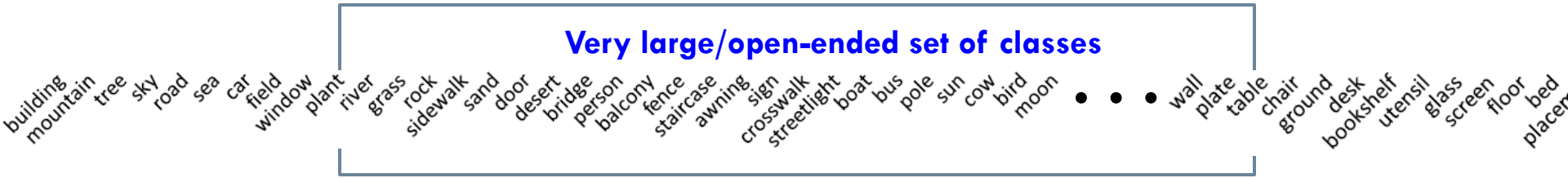
(IMG.XML)

- [car](#)
- [car](#)
- [car](#)
- [car](#)
- [car](#)
- [traffic light](#)
- [traffic light](#)
- [license plate](#)
- [window](#)
- [license plate](#)
- [Street Lamp](#)
- [building](#)
- [buildings](#)
- [road](#)
- [human](#)
- [car](#)
- [window](#)
- [window](#)
- [window](#)
- [windows](#)
- [window](#)
- [window](#)
- [window](#)
- [window](#)
- [window](#)
- [window](#)
- [lamp post](#)
- [lamp post](#)

**Evolving training set**

<http://labelme.csail.mit.edu/>

# Open-universe recognition

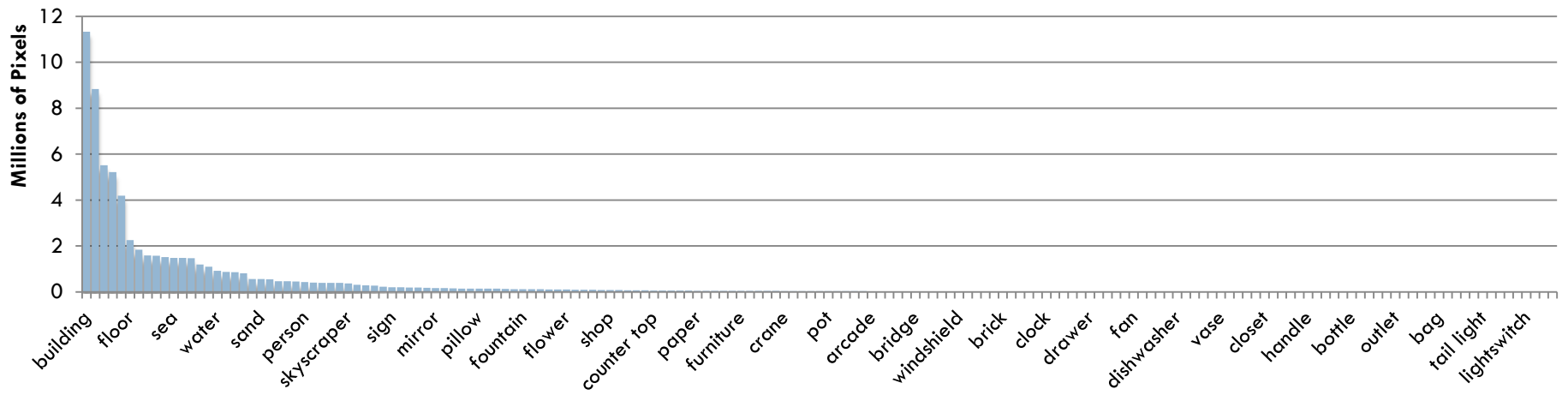


# Open-universe recognition

Very large/open-ended set of classes

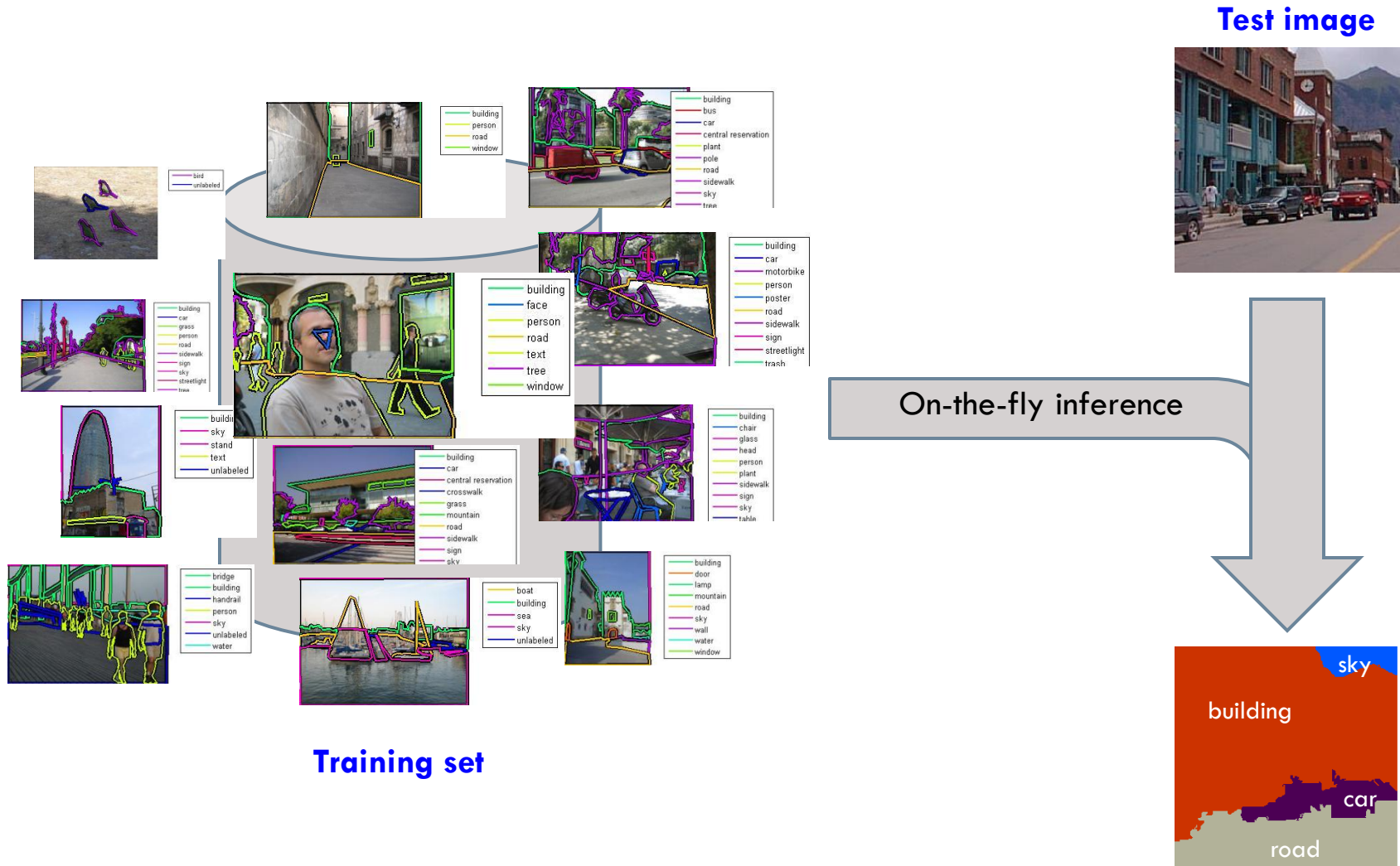
building mountain tree sky road sea car field window plant river grass rock sidewalk sand door desert bridge person balcony fence staircase awning sign crosswalk streetlight boat bus pole sun cow bird moon • • • wall plate table chair ground desk bookshelf utensil glass screen floor bed placemat

Unbalanced data distribution





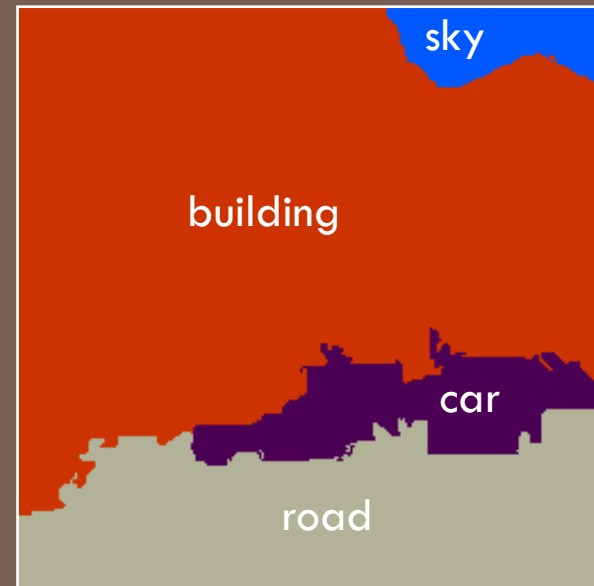
# Potential solution: Lazy learning



# LARGE-SCALE NONPARAMETRIC IMAGE PARSING

Joseph Tighe and Svetlana Lazebnik

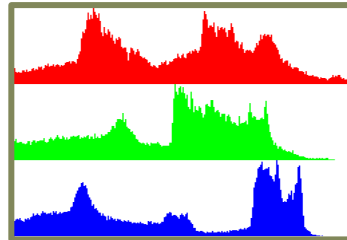
ECCV 2010



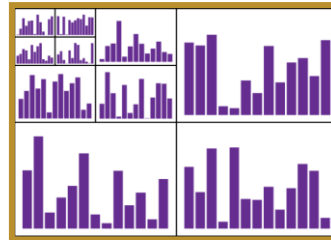
# Step 1: Scene-level matching



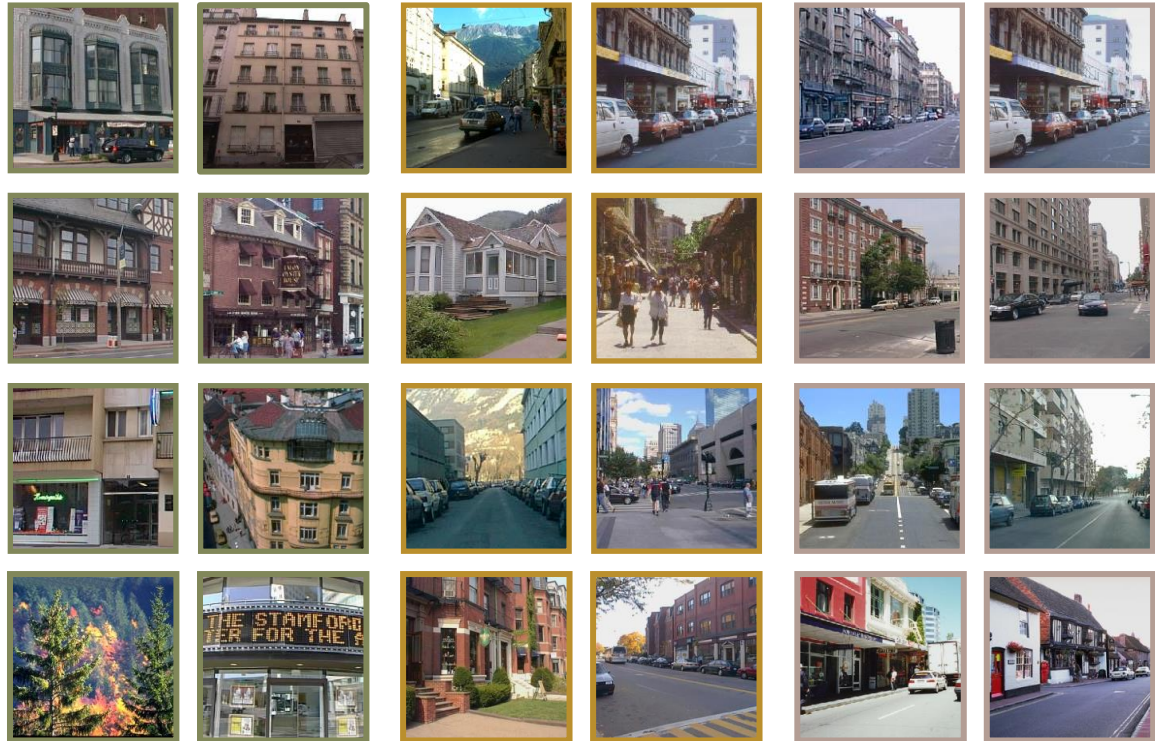
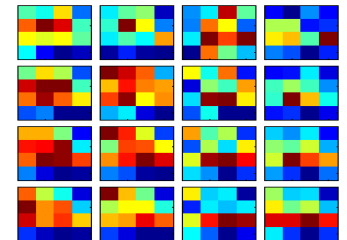
**Color Histogram**



**Spatial Pyramid**  
(Lazebnik et al., 2006)



**Gist**  
(Oliva & Torralba, 2001)



# Step 2: Region-level matching

## Superpixel features



### Superpixels

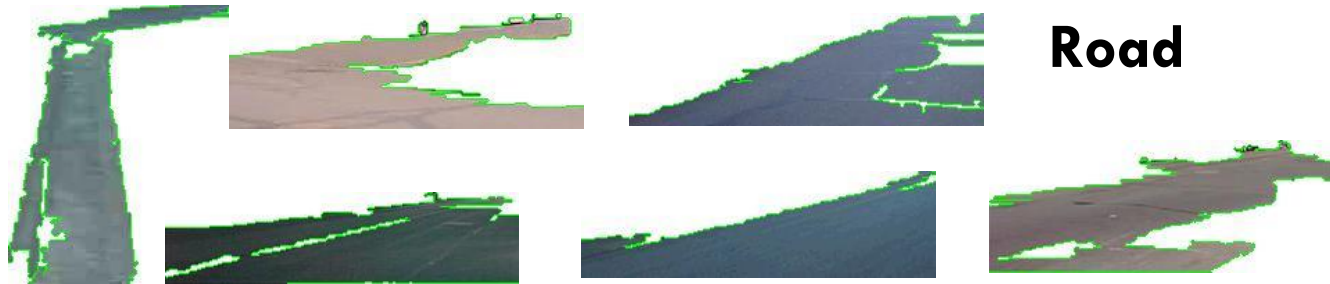
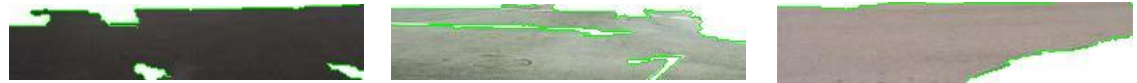
(Felzenszwalb & Huttenlocher, 2004)

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 texton histogram	$100 \times 2$
	SIFT histogram, dilated SIFT histogram	$100 \times 2$
	Left/right/top/bottom boundary SIFT histogram	$100 \times 4$
Color	RGB color mean and std. dev.	$3 \times 2$
	Color histogram (RGB, 11 bins per channel), dilated hist.	$33 \times 2$
Appearance	Color thumbnail ( $8 \times 8$ )	192
	Masked color thumbnail	192
	Grayscale gist over superpixel bounding box	320

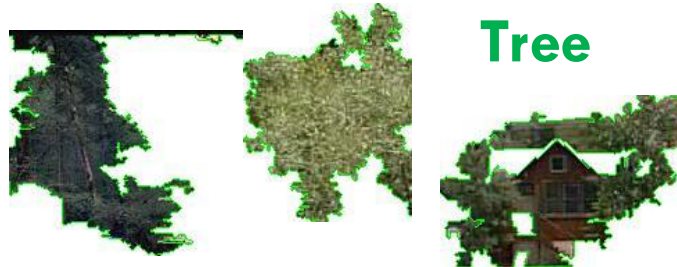
# Step 2: Region-level matching



Pixel Area (size)



Road



Tree



Sky



Building



Snow

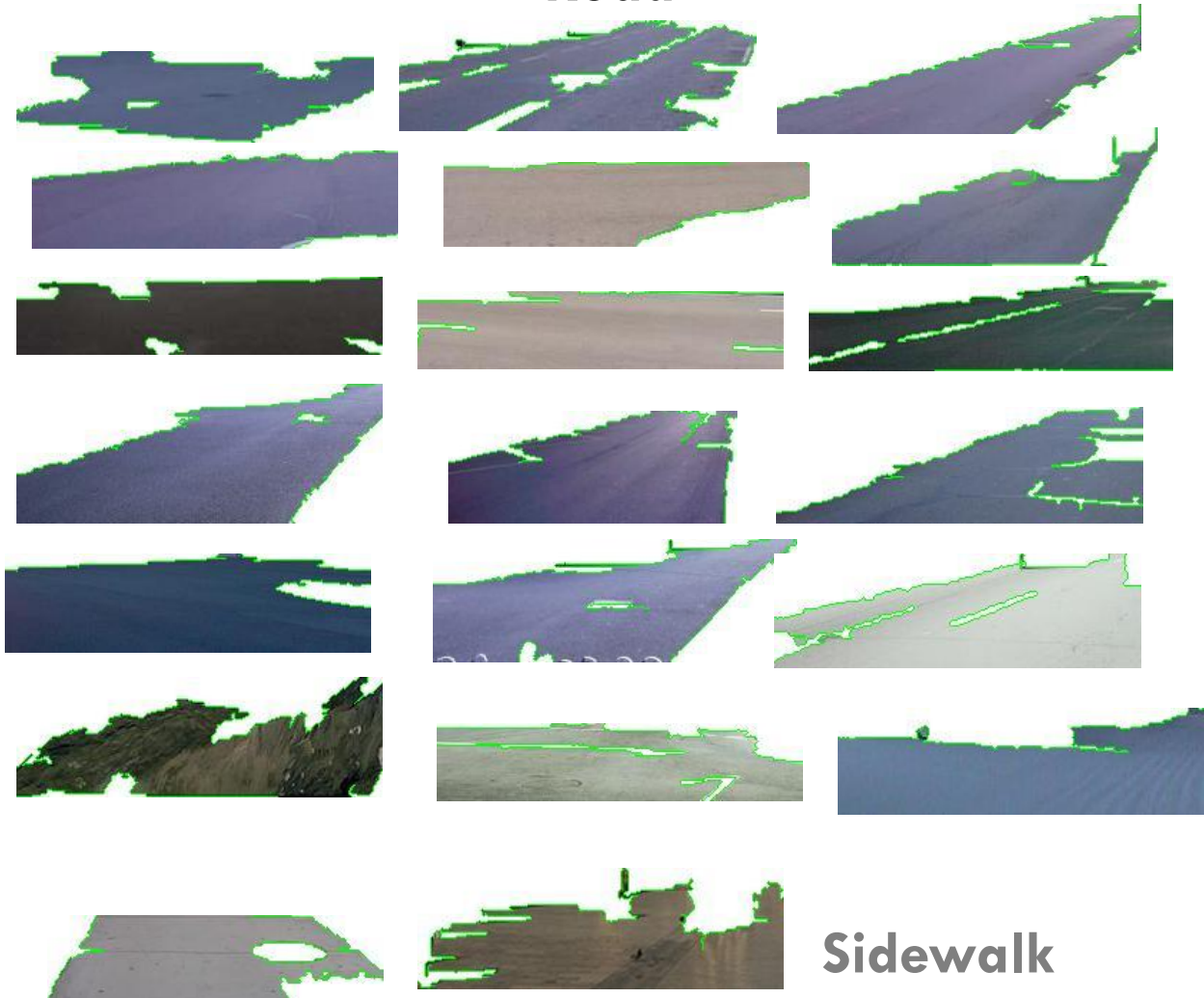
# Step 2: Region-level matching



Absolute mask  
(location)



Road

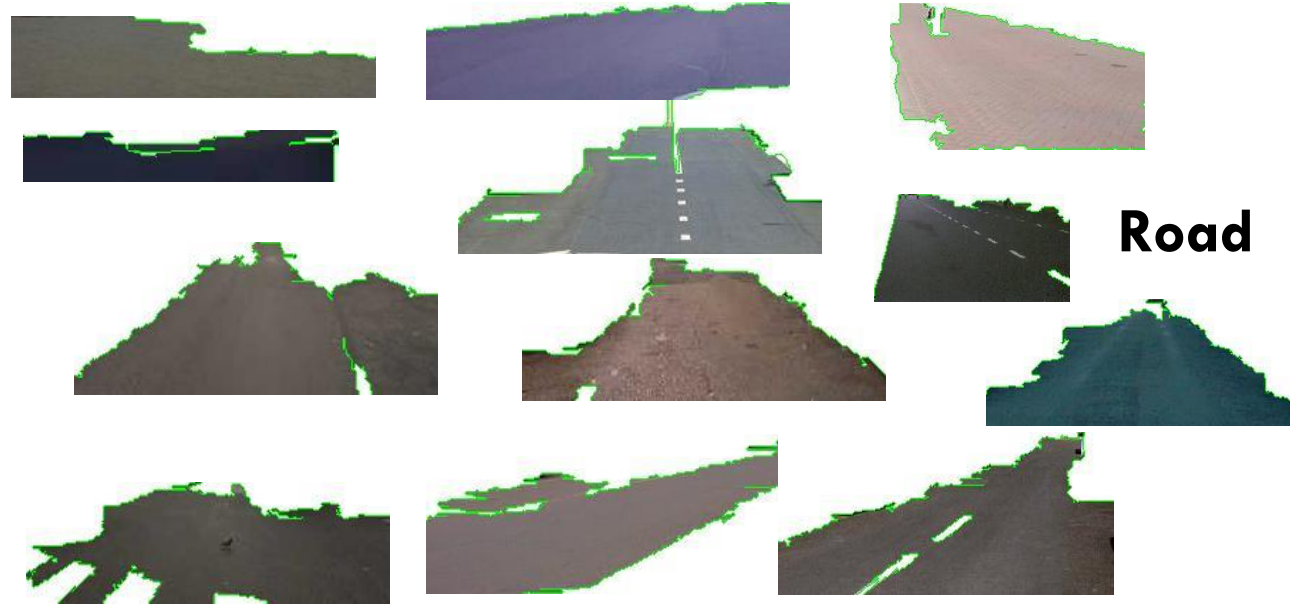


Sidewalk

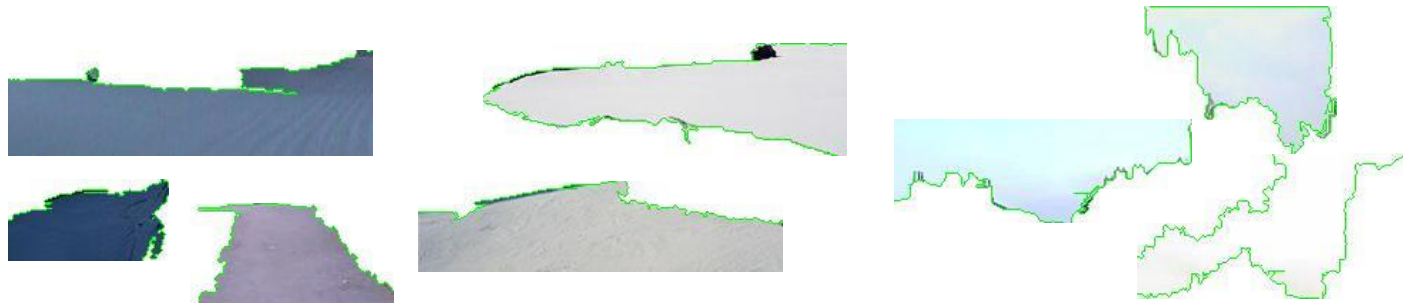
# Step 2: Region-level matching



Texture



Road



Sidewalk

Snow

Sky

# Step 2: Region-level matching



Color histogram

Road



Sidewalk



Building





# Region-level likelihoods

- Nonparametric estimate of class-conditional densities for each class  $c$  and feature type  $k$ :

$$\hat{P}(f_k(r_i) | c) = \frac{\#(N(f_k(r_i)), c)}{\#(D, c)}$$

*k*th feature type of *i*th region

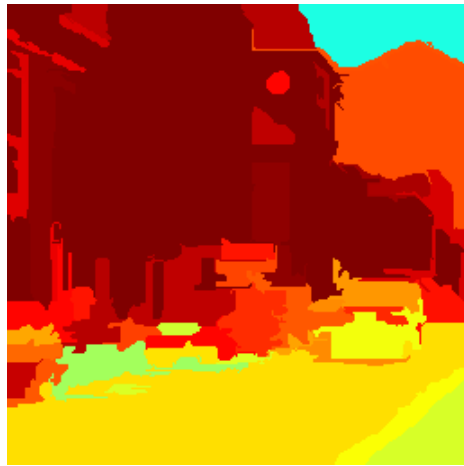
Features of class  $c$  within some radius of  $r_i$

Total features of class  $c$  in the dataset

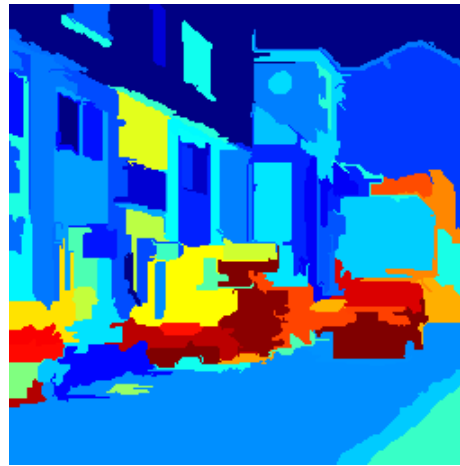
- Per-feature likelihoods combined via Naïve Bayes:

$$\hat{P}(r_i | c) = \prod_{\text{features } k} \hat{P}(f_k(r_i) | c)$$

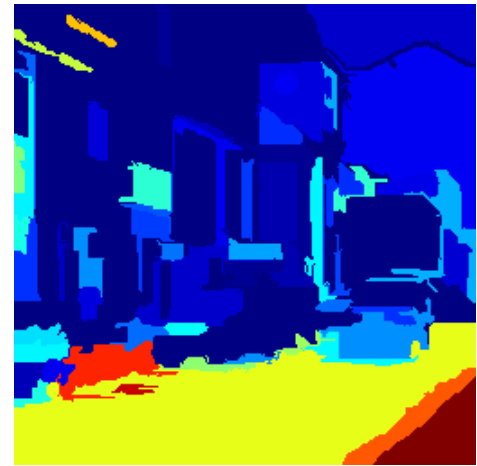
# Region-level likelihoods



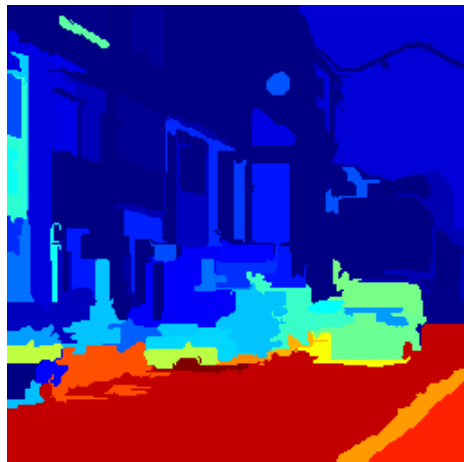
Building



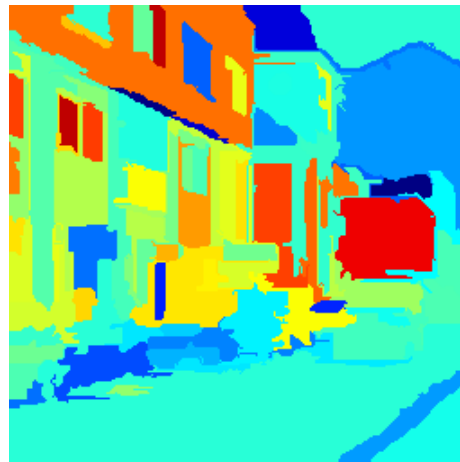
Car



Crosswalk



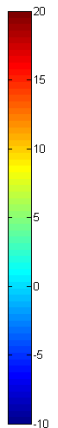
Road



Window



Sky



# Step 3: Global image labeling

- Compute a global image labeling by optimizing a Markov random field (MRF) energy function:

$$E(\mathbf{c}) = \sum_i \underbrace{-\log L(r_i, c_i)}_{\text{Likelihood score for region } r_i \text{ and label } c_i} + \lambda \sum_{i,j} \underbrace{\delta[c_i \neq c_j]}_{\text{Smoothing penalty}} \underbrace{\varphi(c_i, c_j)}_{\text{Co-occurrence penalty}}$$

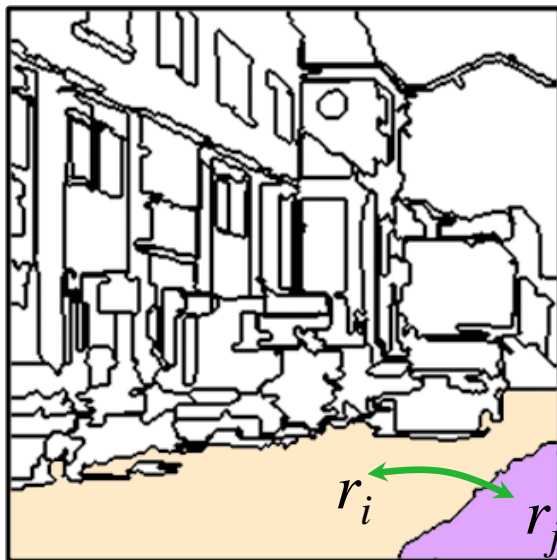
↑  
Vector of region labels

Regions

Neighboring regions

Smoothing penalty

Co-occurrence penalty



Efficient approximate minimization using  $\alpha$ -expansion (Boykov et al., 2002)

# Step 3: Global image labeling

- Compute a global image labeling by optimizing a Markov random field (MRF) energy function:

$$E(\mathbf{c}) = \sum_i \underbrace{-\log L(r_i, c_i)}_{\substack{\text{Likelihood score for} \\ \text{region } r_i \text{ and label } c_i}} + \lambda \sum_{i,j} \underbrace{\delta[c_i \neq c_j]}_{\substack{\text{Smoothing} \\ \text{penalty}}} \underbrace{\varphi(c_i, c_j)}_{\substack{\text{Co-occurrence} \\ \text{penalty}}}$$

↑  
Vector of region labels

Regions

Neighboring regions

# Step 3: Global image labeling

- Compute a global image labeling by optimizing a Markov random field (MRF) energy function:

$$E(\mathbf{c}) = \sum_i \underbrace{-\log L(r_i, c_i)}_{\substack{\text{Likelihood score for} \\ \text{region } r_i \text{ and label } c_i}} + \lambda \sum_{i,j} \underbrace{\delta[c_i \neq c_j]}_{\substack{\text{Smoothing} \\ \text{penalty}}} \underbrace{\varphi(c_i, c_j)}_{\substack{\text{Co-occurrence} \\ \text{penalty}}}$$

↑  
Vector of region labels

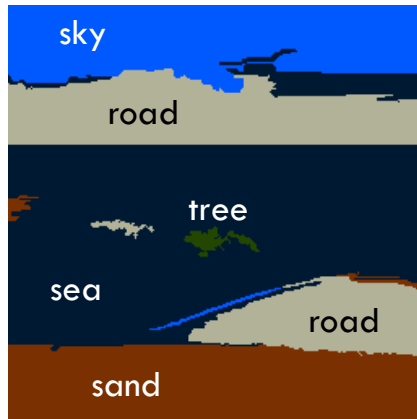
Regions

Neighboring regions

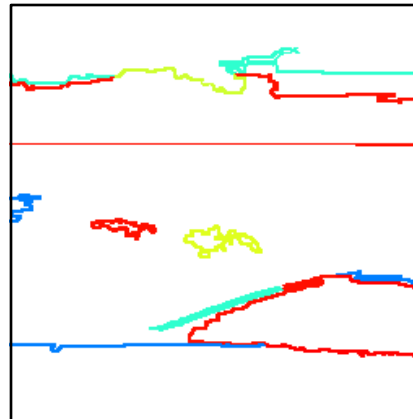
Original image



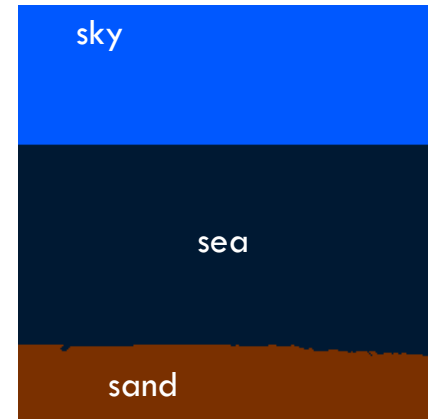
Maximum likelihood labeling



Edge penalties

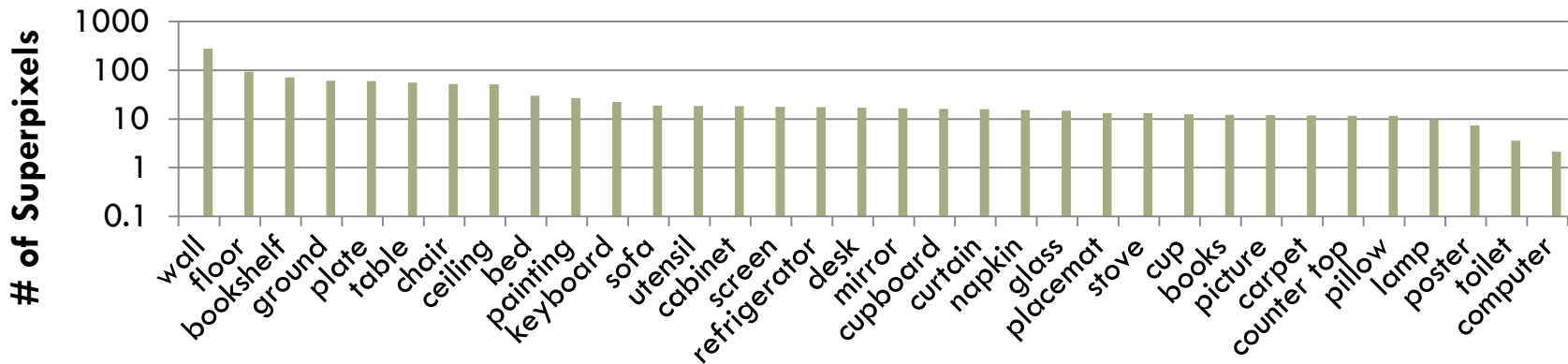
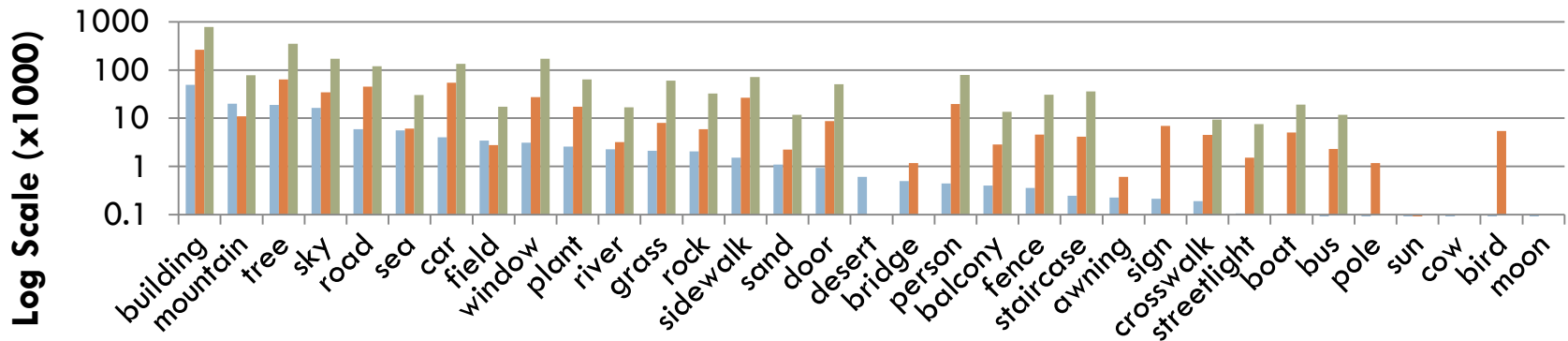


MRF labeling

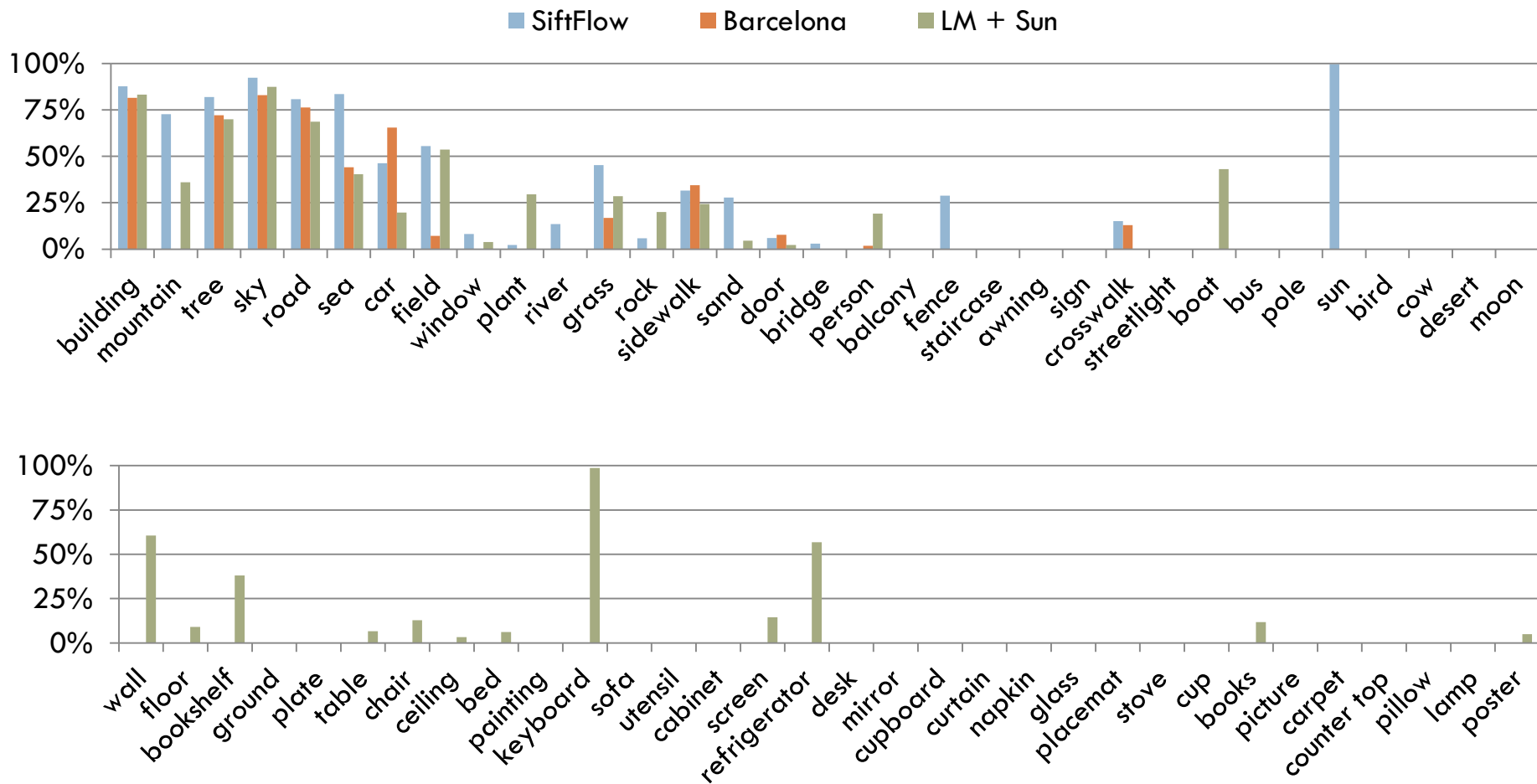


# Datasets

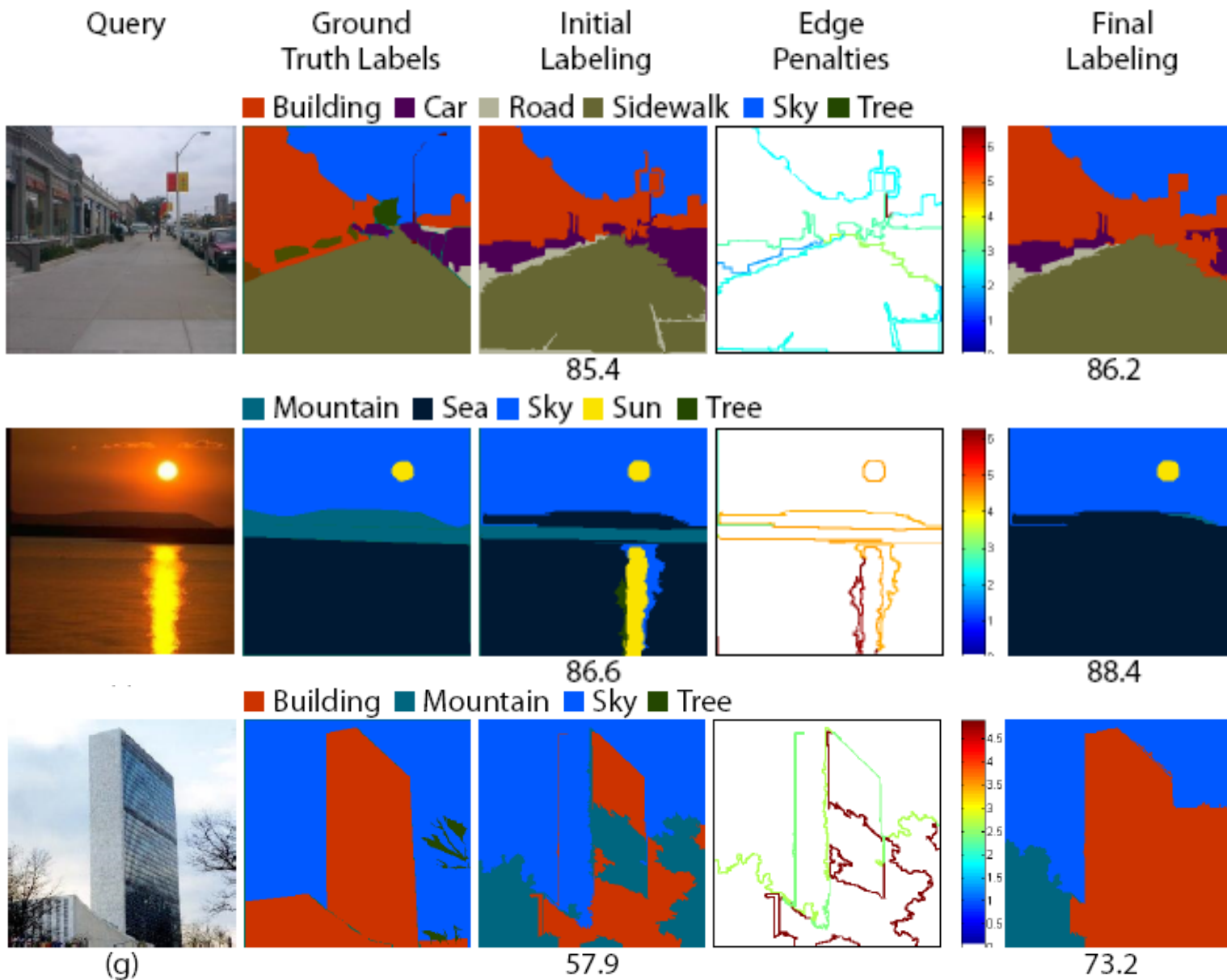
	Training images	Test images	Labels
<b>SIFT Flow</b> (Liu et al., 2009)	2,488	200	33
<b>Barcelona</b>	14,871	279	170
<b>LabelMe+SUN</b>	50,424	300	232



# Per-class classification rates

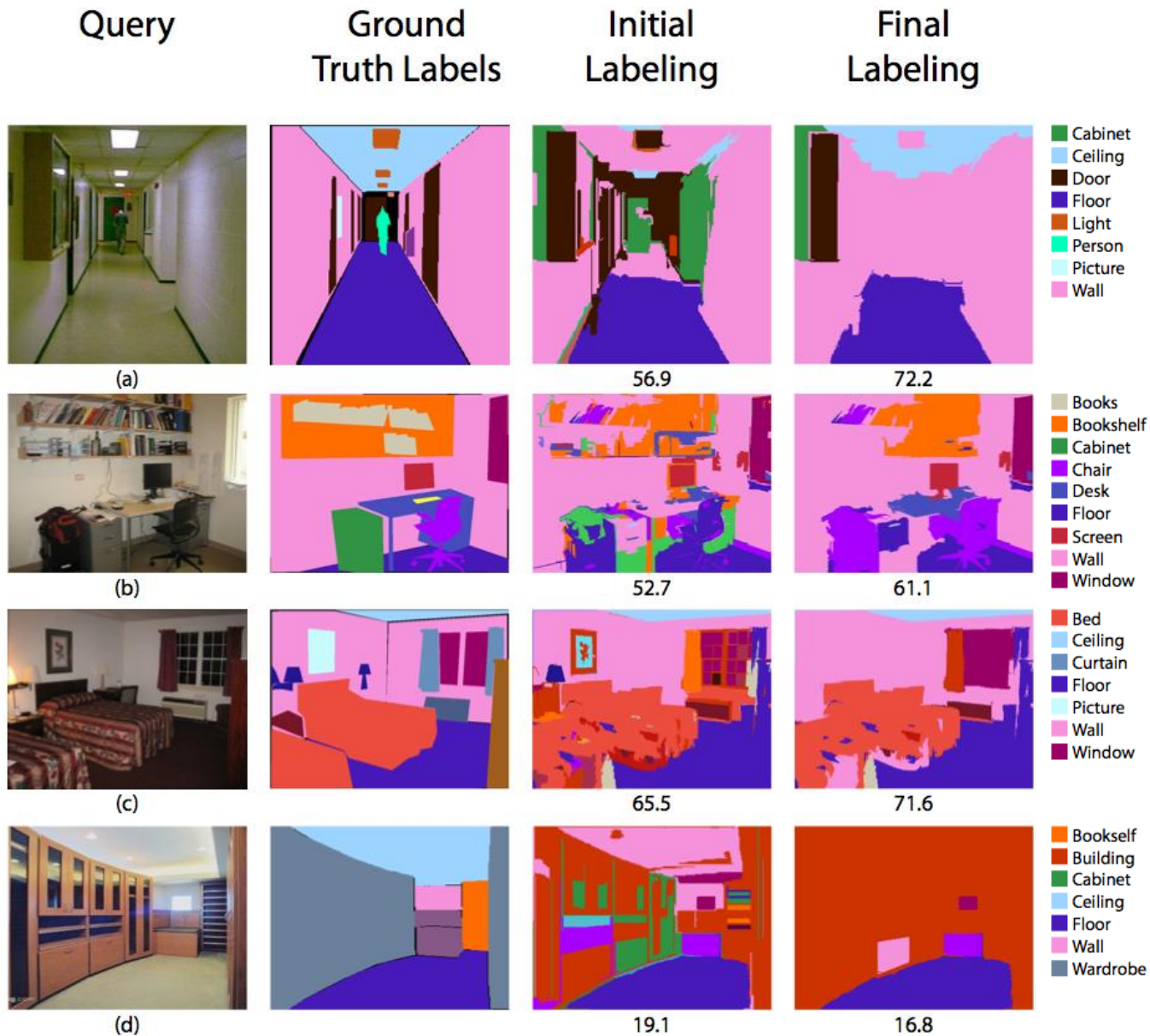


# Results on SIFT Flow dataset





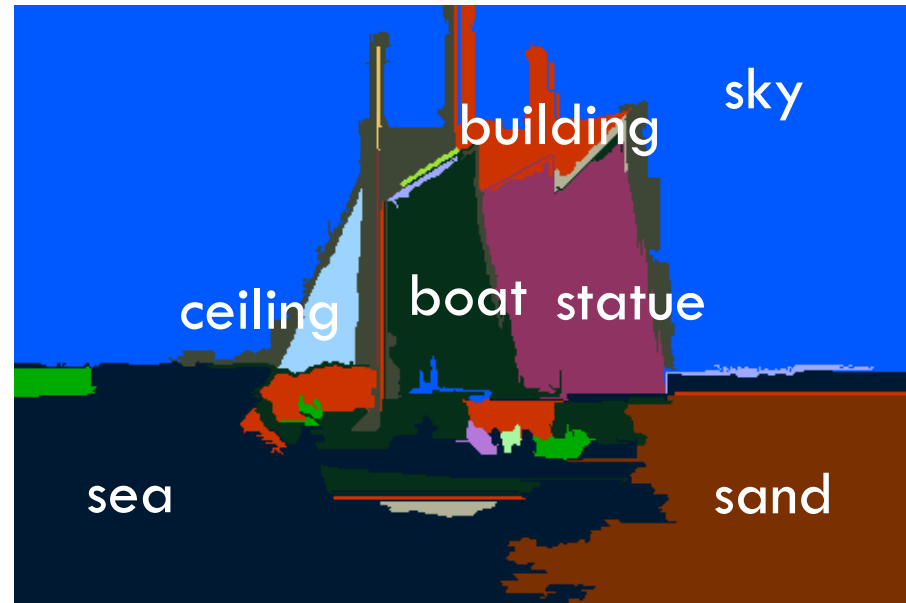
# Results on LM+SUN dataset



# Summary so far

- A lazy learning method for image parsing:
  - Global scene matching
  - Superpixel-level matching
  - MRF optimization
- Challenges
  - Indoor images are hard!
  - We do well on “stuff” but not on “things”

# We get the “stuff” but not the “things”



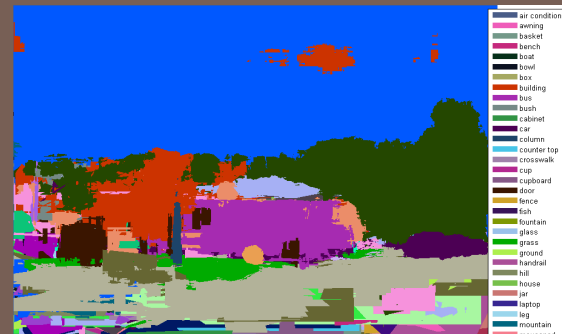
# FINDING THINGS: IMAGE PARSING WITH REGIONS AND PER-EXEMPLAR DETECTORS

Joseph Tighe and Svetlana Lazebnik  
CVPR 2013

Superparsing Result

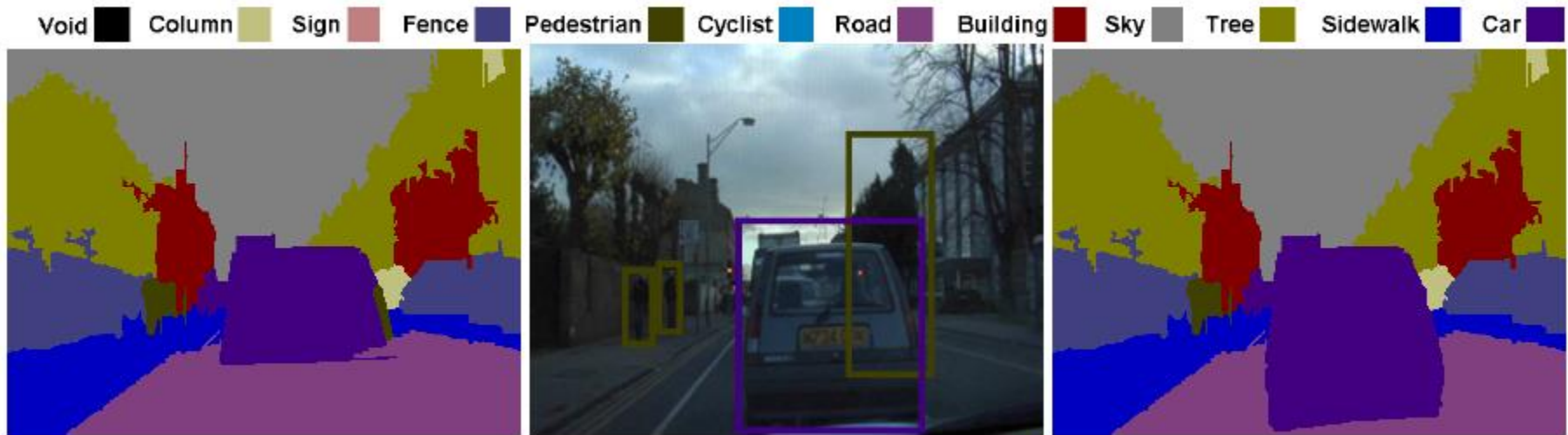


Detector Based Parsing Result



# To get the “things” use detectors

- Ladicky et al. used detector output coupled with bounding box based foreground/background segmentation to improve performance on things



Result without  
detections

Set of detections

Final Result

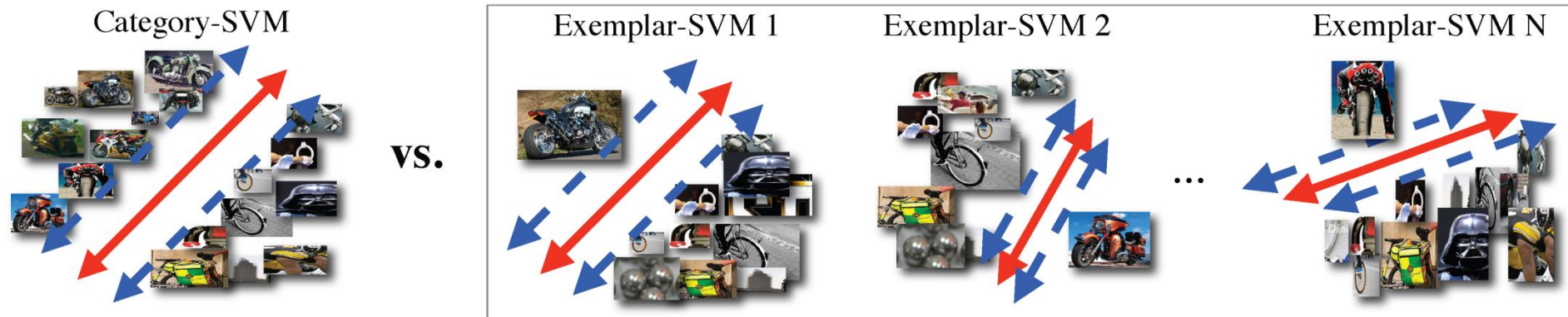
# Problems with this approach

- The mask for bounding boxes is obtained by an automatic segmentation, which can fail
- The models must be pre-trained and cannot adapt to new data easily
- There is little flexibility for objects that take many forms

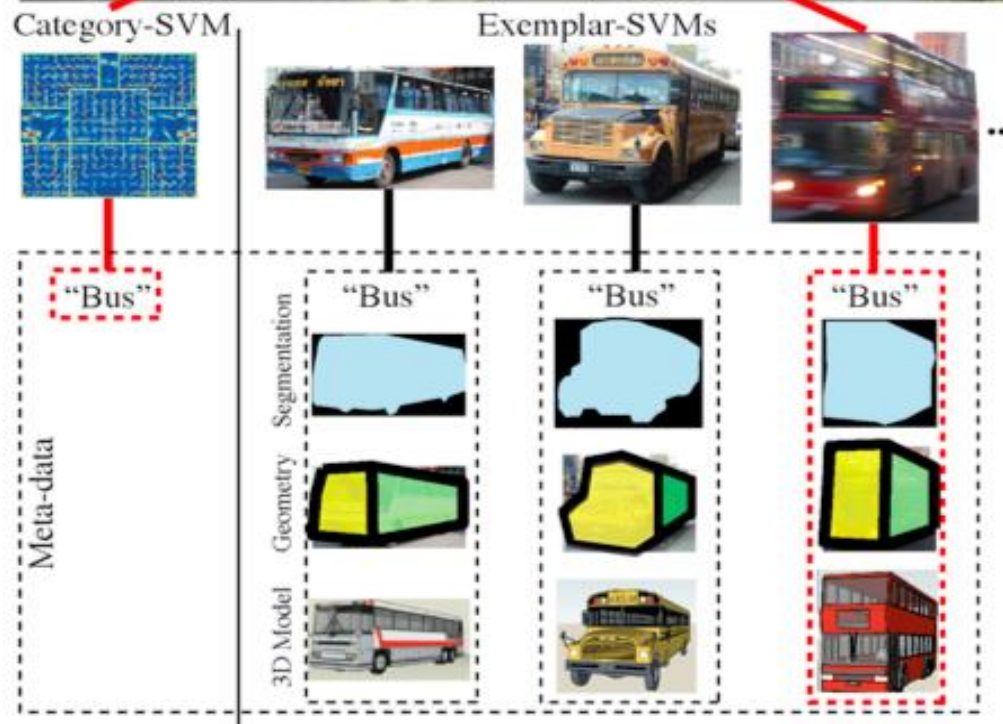
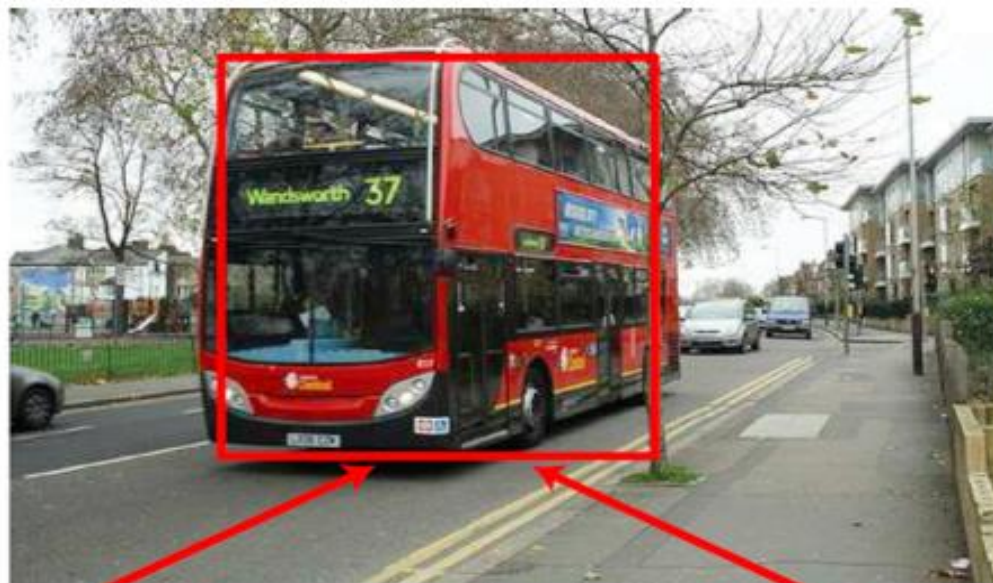


# Per-exemplar detectors

- For each instance of a class: train SVM based on HOG features
- Negative examples are taken from all images that do not contain the class



Tomasz Malisiewicz, Abhinav Gupta, Alexei A. Efros. Ensemble of Exemplar-SVMs for Object Detection and Beyond. In ICCV, 2011



Tomasz Malisiewicz, Abhinav Gupta, Alexei A. Efros. Ensemble of Exemplar-SVMs for Object Detection and Beyond . In ICCV, 2011



# Per-exemplar detectors for parsing

- Retrieve a set of similar images using global image descriptors
- Train per-exemplar detectors for “things” in retrieval set
- Run trained detectors on query and transfer weighted mask for all positive detections

# Retrieval set for



- car
- car
- car
- car

1



- sign
- building
- pole
- window
- car
- streetlight
- road
- building
- sky
- sidewalk
- grass
- sidewalk
- sidewalk
- flower
- flower
- pole
- pole
- pole
- truck
- car
- window
- balcony
- window
- window
- window

2



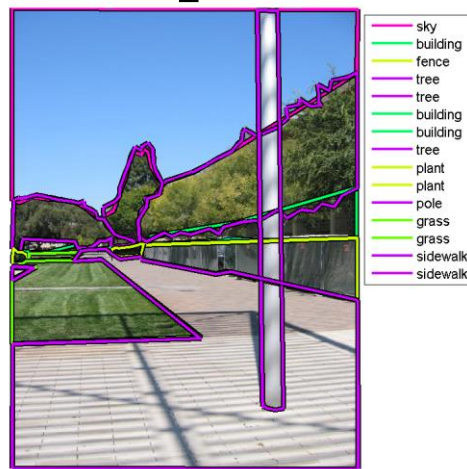
- road
- mountain
- snow
- cloud
- cloud
- sky

3



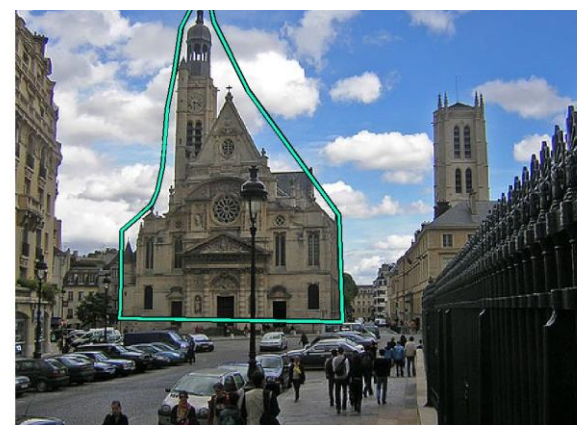
- boat
- building
- boat
- sky
- building
- tree
- tree
- tree
- tree
- dock
- building
- building
- building
- building
- tree
- river
- boat
- dock
- wall
- sidewalk
- person
- person
- person
- person
- person
- person
- person
- person
- sign
- window

4



- sky
- building
- fence
- tree
- tree
- building
- building
- plant
- pole
- grass
- grass
- sidewalk
- sidewalk

5



6



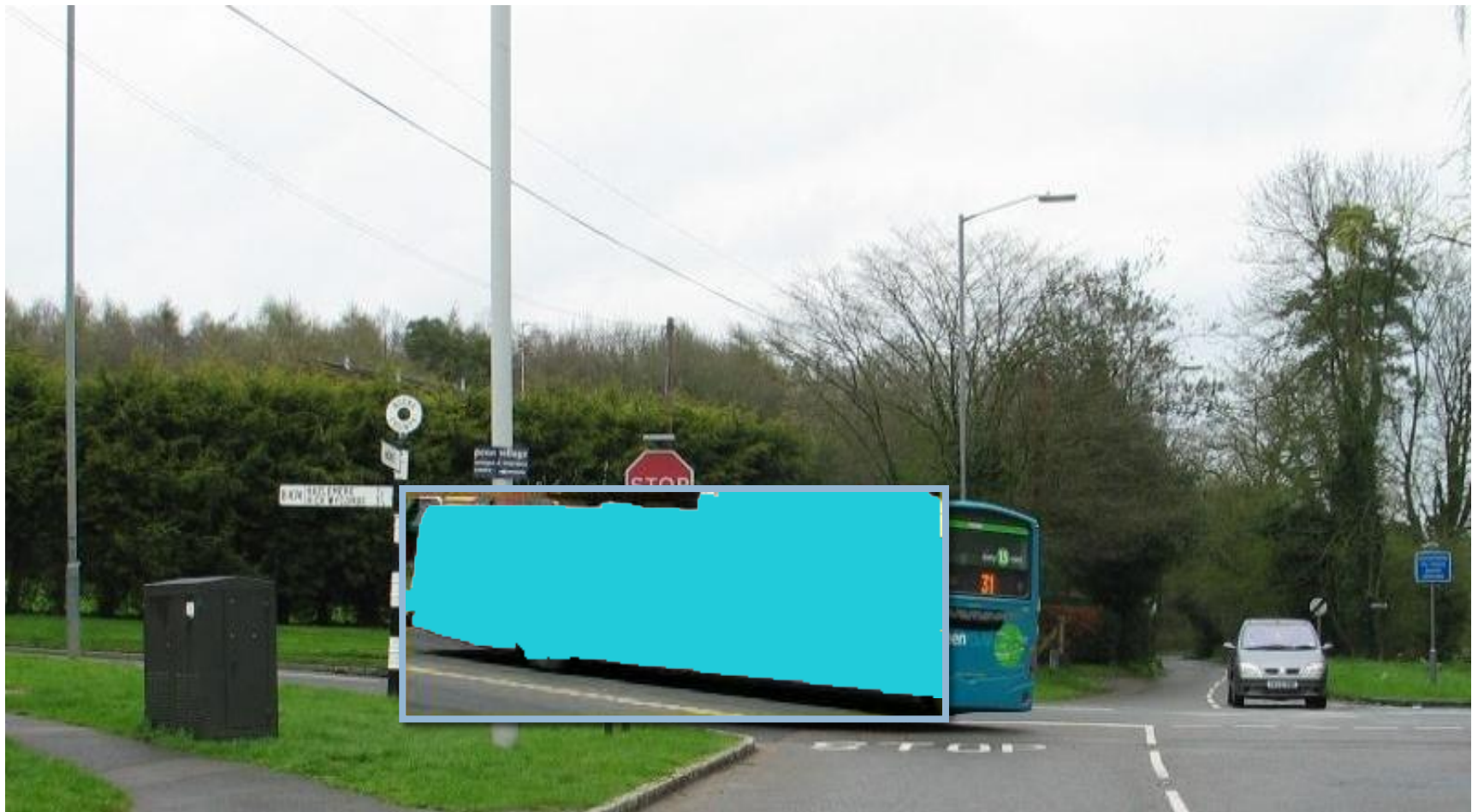
# Per-exemplar detectors for parsing

- Retrieve a set of similar images using global image descriptors
- Train per-exemplar detectors for each object in retrieval set
- Run trained detectors on query and transfer weighted masks for all positive detections

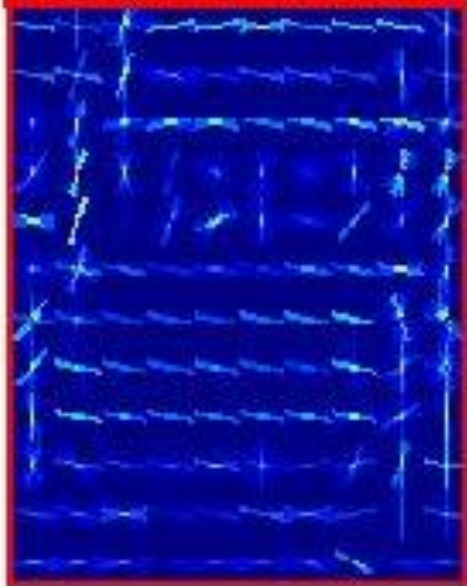
# Per-exemplar detectors for parsing



# Per-exemplar detectors for parsing



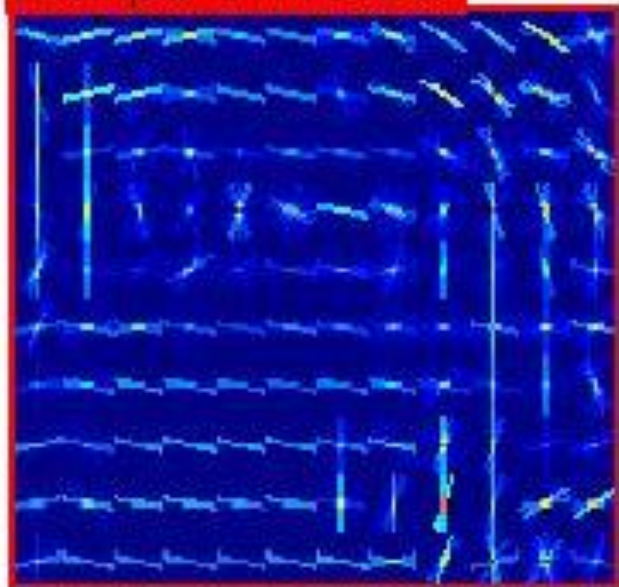
Exemplar-SVM bus 20



Exemplar Image 20



Exemplar-SVM bus 28



Exemplar Image 28



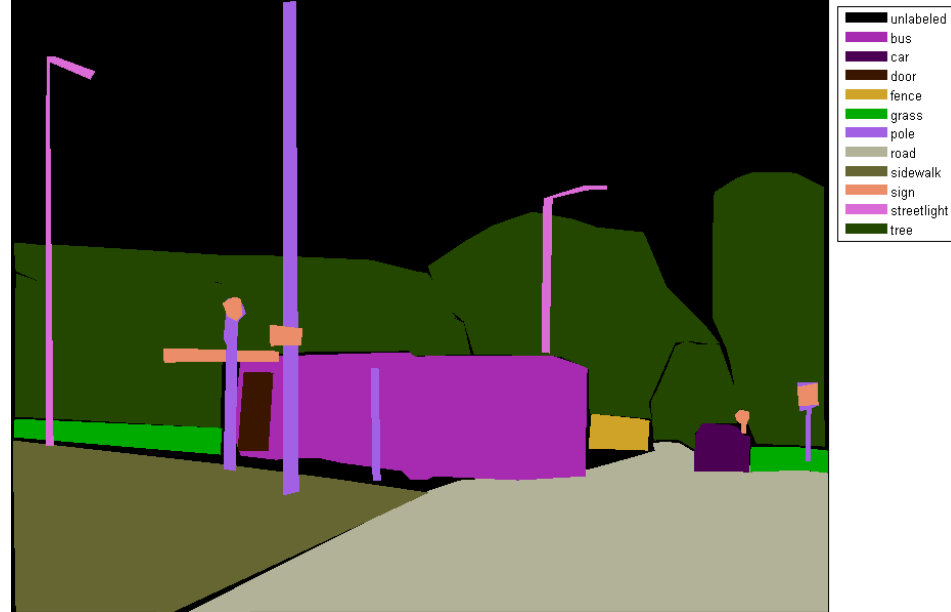


# Per-exemplar detectors for parsing

- Retrieve a set of similar images using global image descriptors
- Train per-exemplar detectors for “things” in retrieval set
- Run trained detectors on query and transfer weighted masks for all positive detections

# Per-exemplar detectors for parsing





Superparsing Result

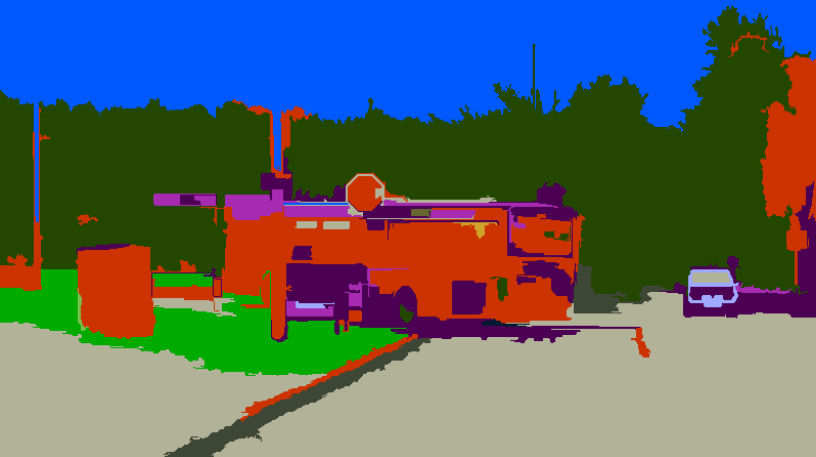
Detector-based Parsing Result



55% (23%)

45% (26%)

# Superparsing Result



- building
- bus
- car
- church
- fence
- grass
- house
- road
- sea
- sidewalk
- sky
- snow
- tree

55% (23%)



# Detector Based Parsing Result



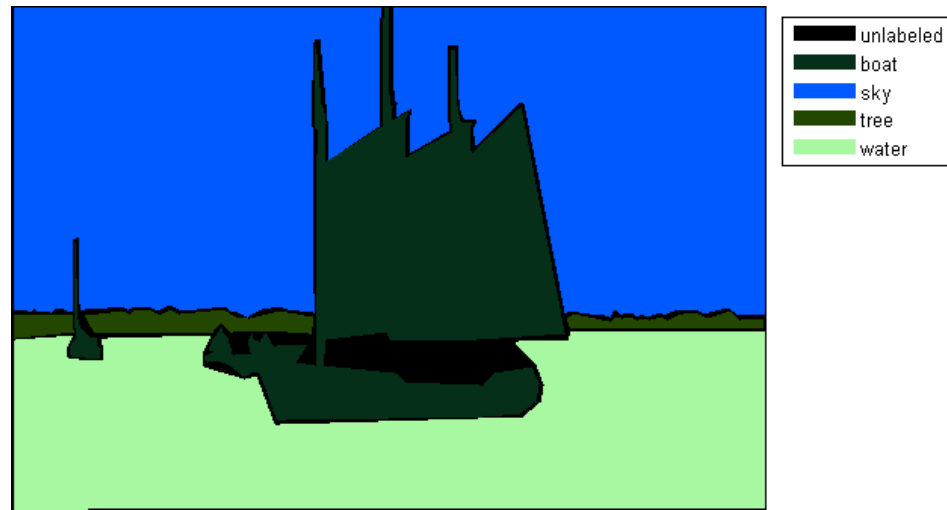
- air conditioner
- awning
- basket
- bench
- boat
- bowl
- box
- building
- bus
- bush
- cabinet
- car
- column
- counter top
- crosswalk
- cup
- cupboard
- door
- fence
- fish
- fountain
- glass
- grass
- ground
- handrail
- hill
- house
- jar
- laptop
- leg
- mountain
- mousepad

45% (26%)



- building
- bus
- car
- church
- column
- door
- fence
- grass
- house
- person
- plant
- pole
- road
- sidewalk
- sign
- sky
- tree
- wheel

61% (31%)



Superparsing Result

Detector Based Parsing Result



- animal
- boat
- bridge
- building
- ceiling
- church
- fruit
- grass
- road
- sand
- sea
- sky
- snow
- statue
- tower
- water



- air conditioner
- airplane
- boat
- books
- bookshelf
- bridge
- building
- car
- ceiling
- door
- field
- grass
- ground
- hill
- mountain
- pen
- plate

52% (31%)

19% (25%)

# Superparsing Result



- animal
- boat
- bridge
- building
- ceiling
- church
- fruit
- grass
- road
- sand
- sea
- sky
- snow
- statue
- tower
- water

52% (31%)

# Detector Based Parsing Result



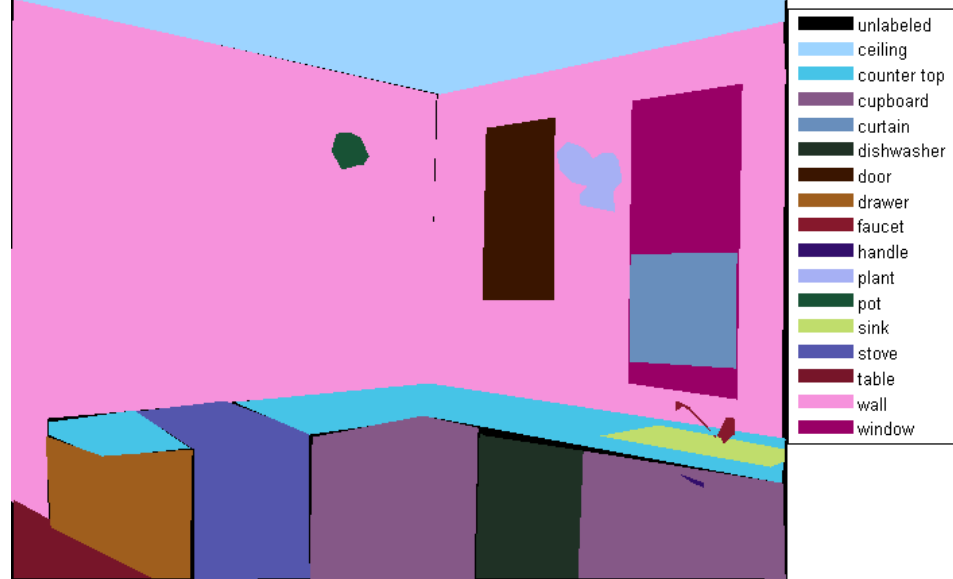
- air conditioner
- airplane
- boat
- books
- bookshelf
- bridge
- building
- car
- ceiling
- door
- field
- grass
- ground
- hill
- mountain
- pen
- plate

19% (25%)

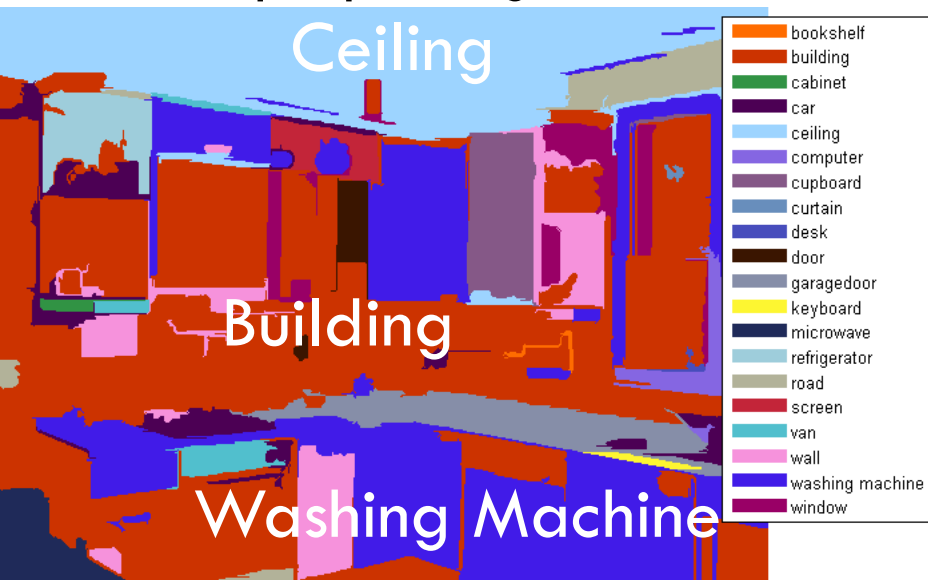


- boat
- building
- church
- grass
- mountain
- road
- sand
- sea
- sky
- wall

62% (46%)



Superparsing Result



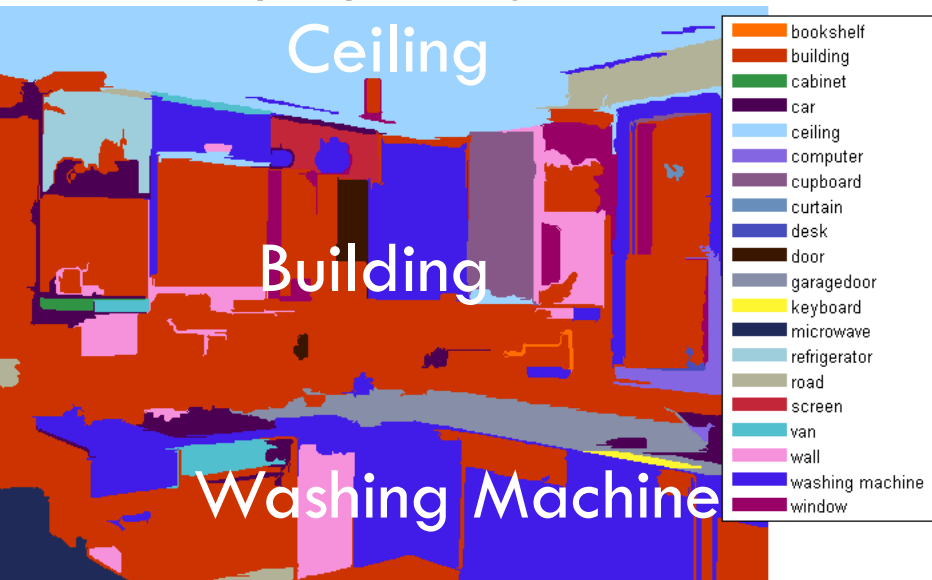
12% (7%)

Detector Based Parsing Result



20% (9%)  
Dishwasher

# Superparsing Result

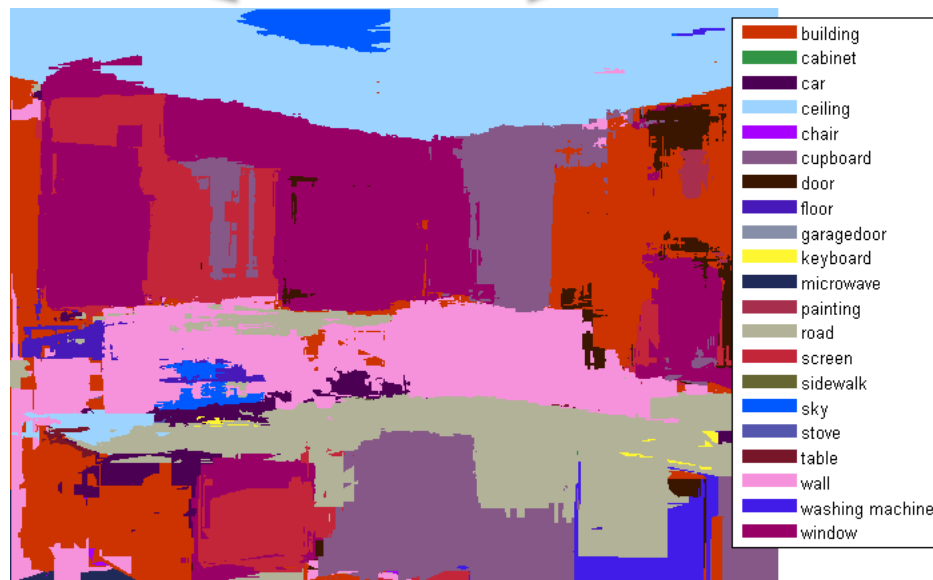


12% (7%)

# Detector Based Parsing Result



20% (9%)



24% (10%)



# Conclusion

---

- Image parsing with superpixels
  - ▣ Scene-level matching
  - ▣ Superpixel-level matching
  - ▣ MRF optimization
- Getting “things” with detectors
  - ▣ Use per-exemplar detectors of Malisiewicz et al.