

Recognition Methods for Open-Universe Datasets

Lana Lazebnik

Closed-universe recognition

Fixed, pre-defined set of classes

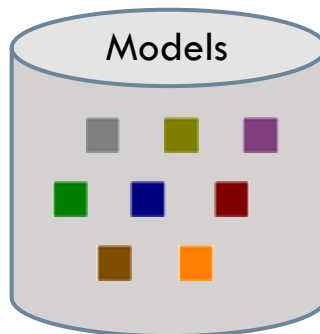
■ sky ■ tree ■ road ■ grass ■ water ■ bldg ■ mntn ■ fg obj.

**Fixed, static
training set**



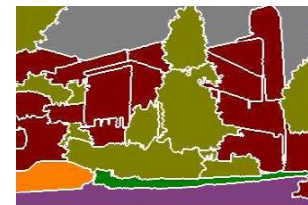
Learning
(offline)

Models



Inference

Test image



Output

Closed-universe datasets



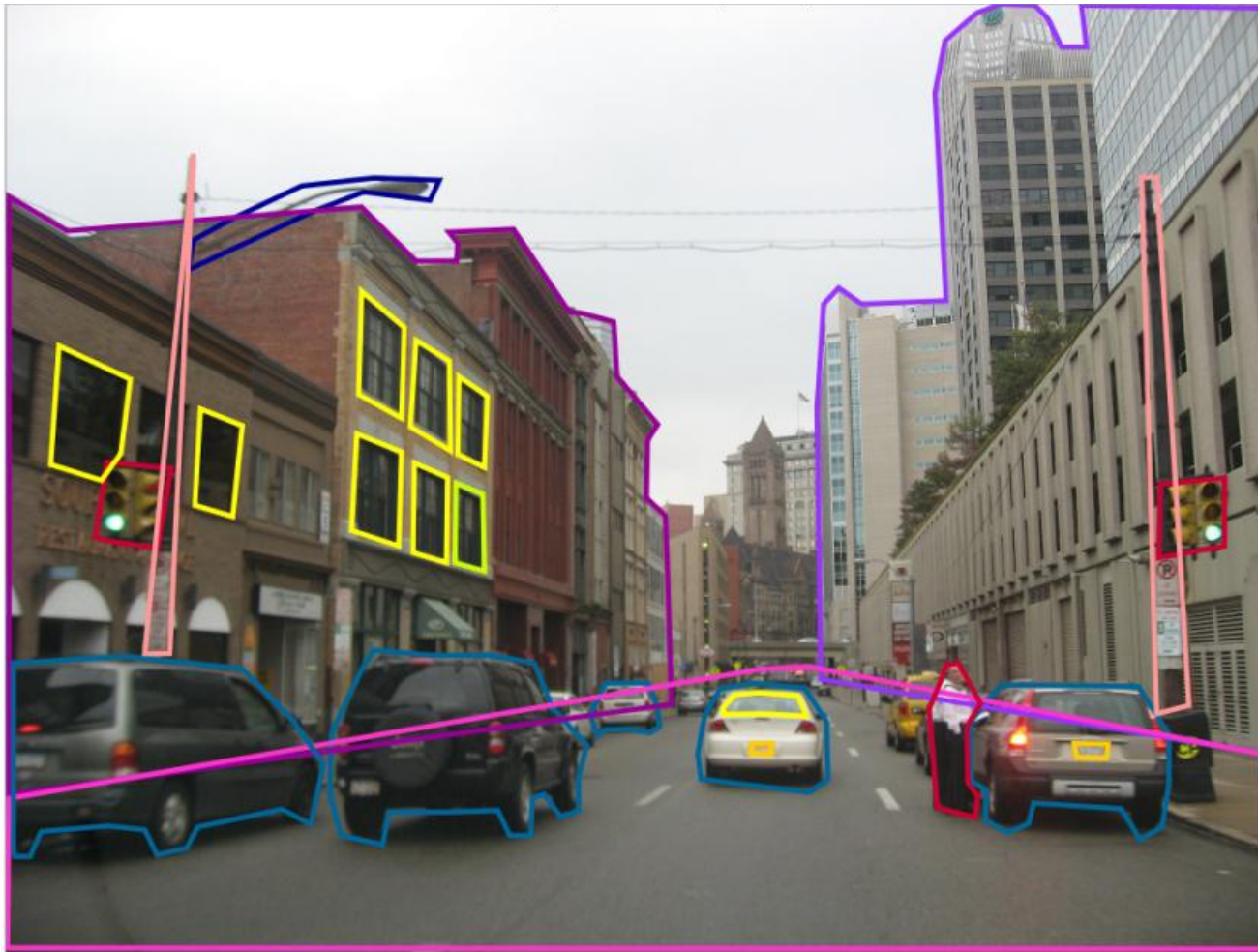
- Small amount of data
- Static datasets
- Limited variation
- Full annotation

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

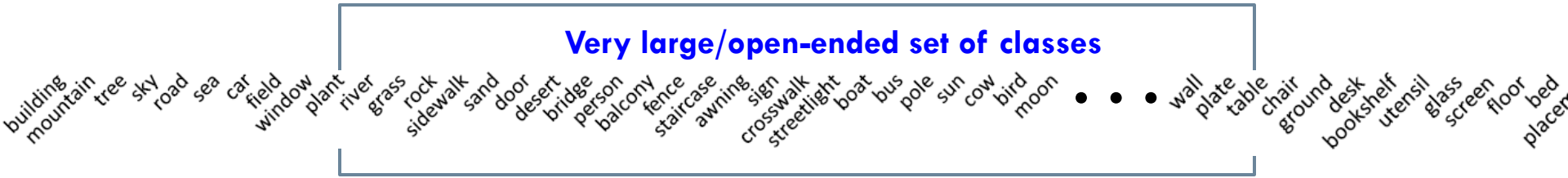
(IMG.XML)

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

Evolving training set

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

Open-universe recognition

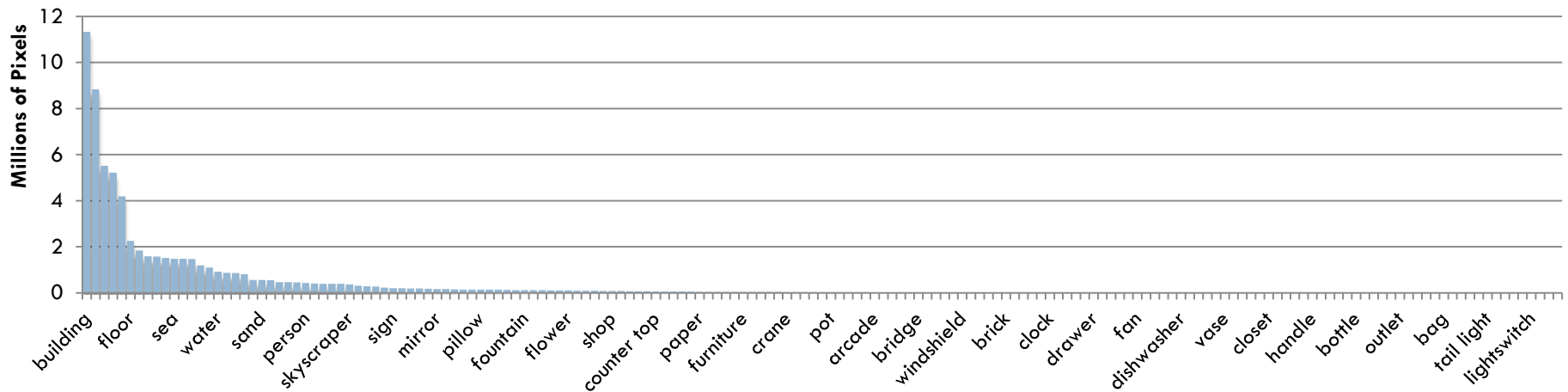


Open-universe recognition

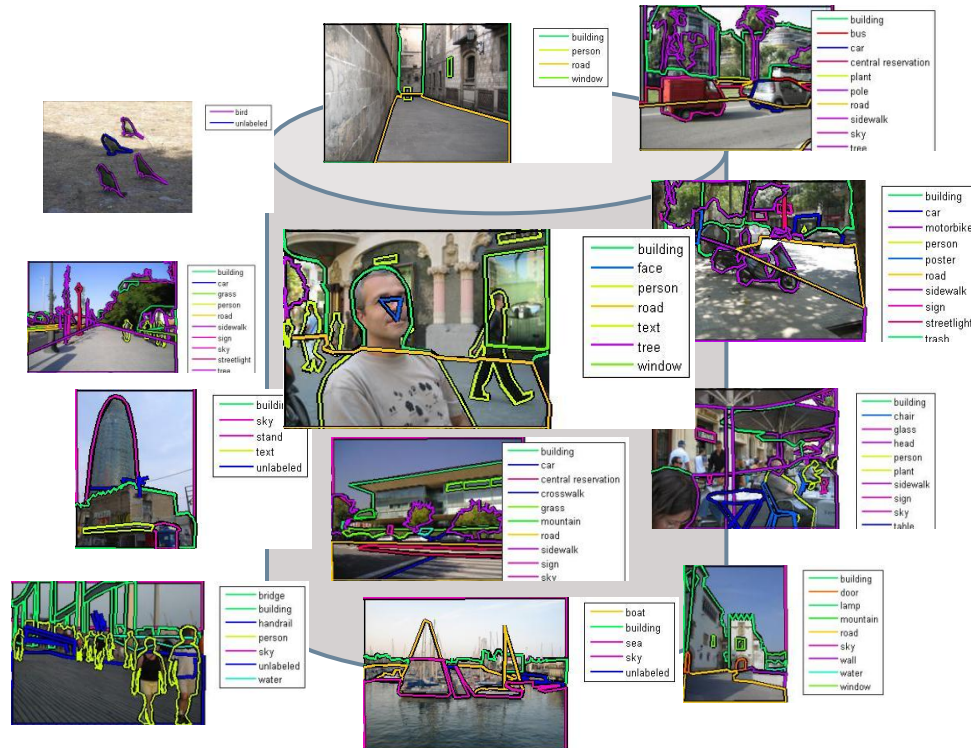
Very large/open-ended set of classes

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

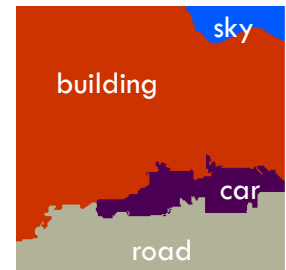


Training set

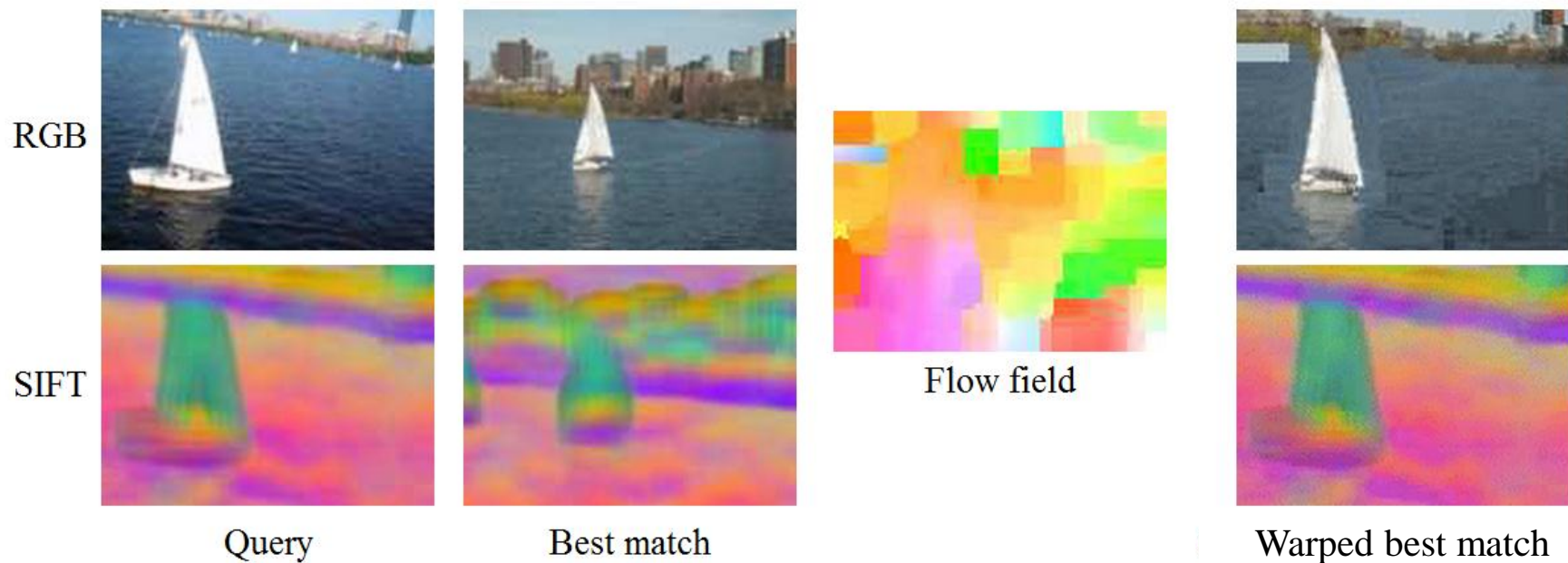
Test image



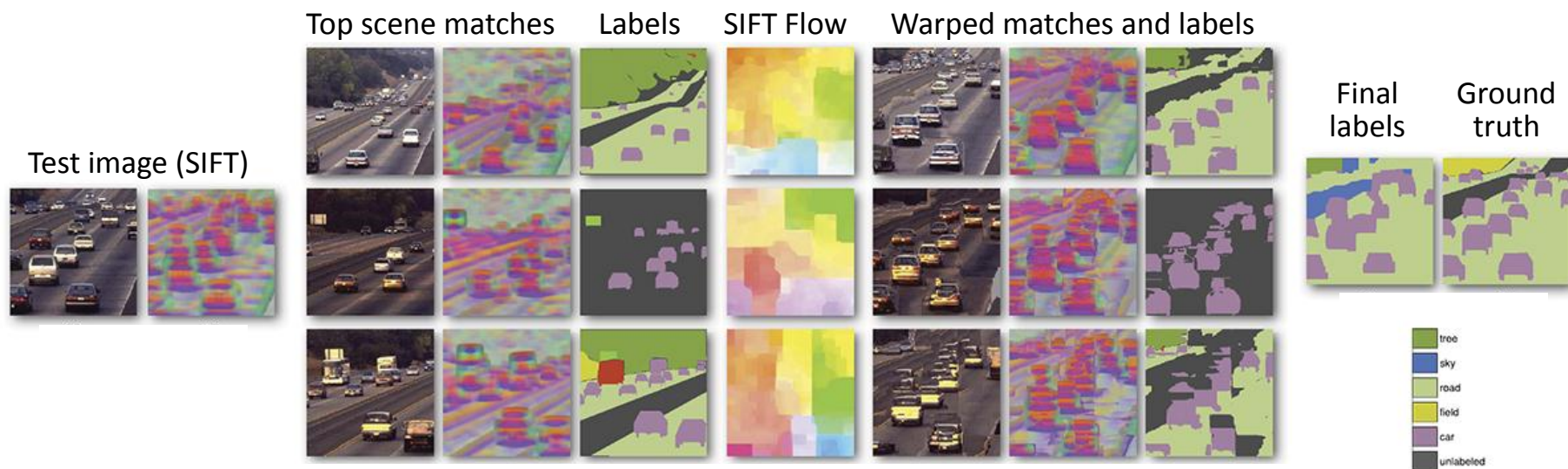
On-the-fly inference



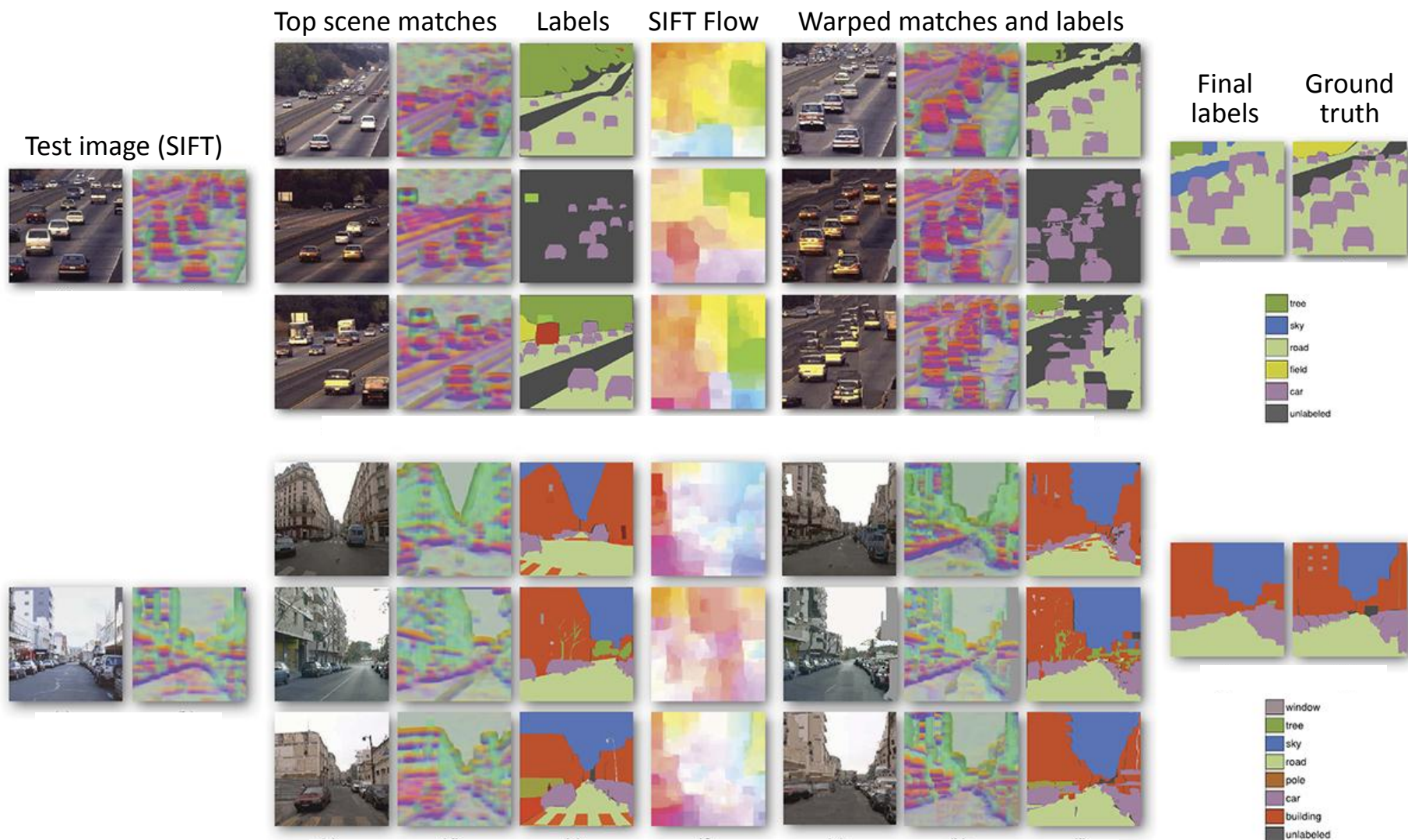
Dense Scene Alignment by SIFT Flow



Dense Scene Alignment by SIFT Flow



Dense Scene Alignment by SIFT Flow



SIFT Flow: Pros and cons

□ Advantages

- ▣ Nonparametric method, can work with any number of labels, evolving training set
- ▣ Initial global scene matching step improves efficiency, provides scene-level context

□ Disadvantages

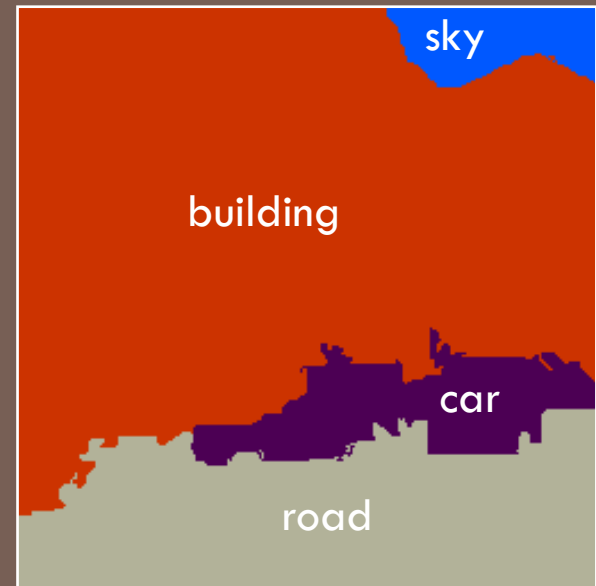
- ▣ Computing SIFT flow is very computationally intensive
- ▣ Warping model is not necessarily the most natural one for describing object-level correspondence

LARGE-SCALE NONPARAMETRIC IMAGE PARSING



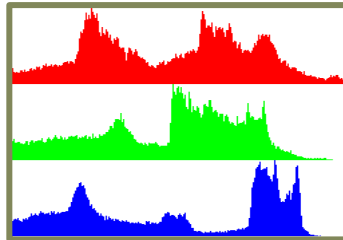
Joseph Tighe and Svetlana Lazebnik

ECCV 2010, new stuff in preparation

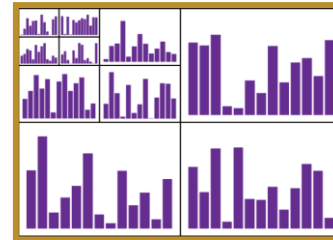


Step 1: Scene-level matching

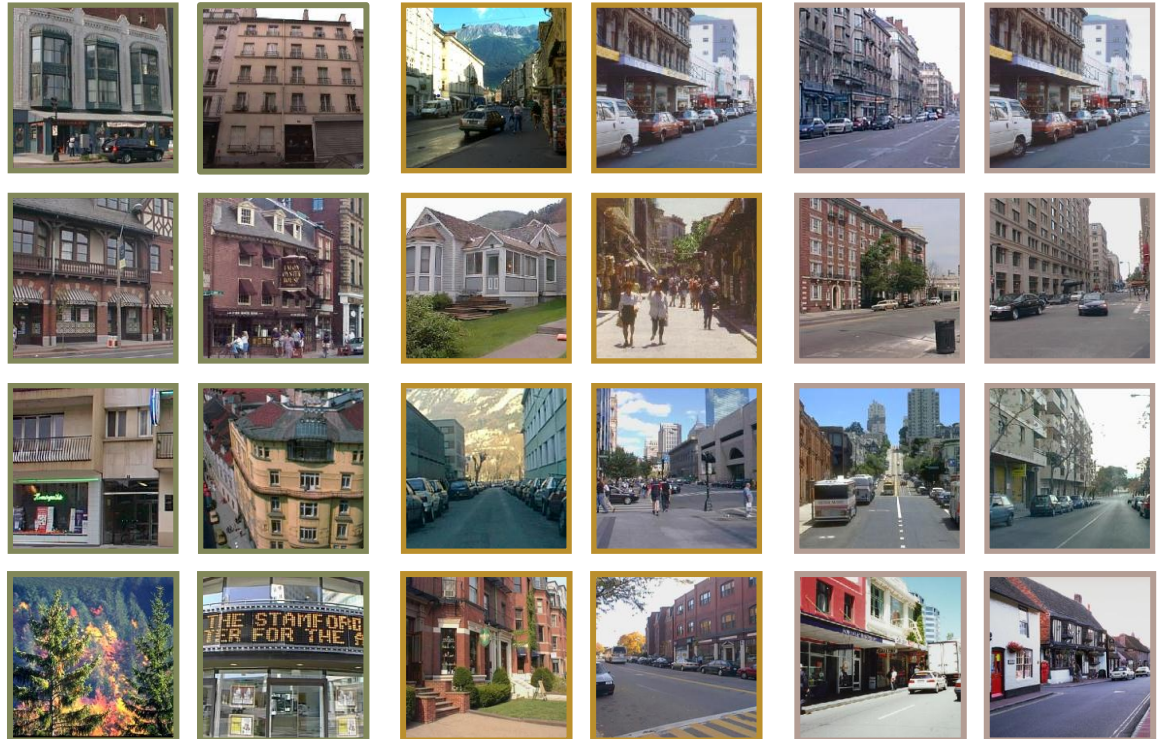
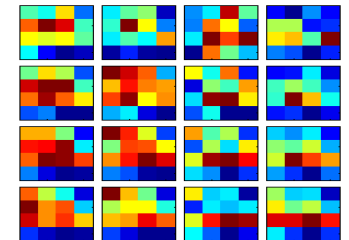
Color Histogram



Spatial Pyramid
(Lazebnik et al., 2006)



Gist
(Oliva & Torralba, 2001)



Step 2: Region-level matching

Supersixel features



Supersixels

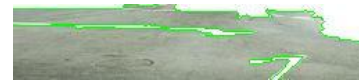
(Felzenszwalb & Huttenlocher, 2004)

Shape	Mask of supersixel shape over its bounding box (8×8)	64
	Bounding box width/height relative to image width/height	2
	Supersixel area relative to the area of the image	1
Location	Mask of supersixel shape over the image	64
	Top height of bounding box relative to image height	1
Texture/SIFT	Texton histogram, dilated texton histogram	100×2
	SIFT histogram, dilated SIFT histogram	100×2
	Left/right/top/bottom boundary SIFT histogram	100×4
Color	RGB color mean and std. dev.	3×2
	Color histogram (RGB, 11 bins per channel), dilated hist.	33×2
Appearance	Color thumbnail (8×8)	192
	Masked color thumbnail	192
	Grayscale gist over supersixel bounding box	320

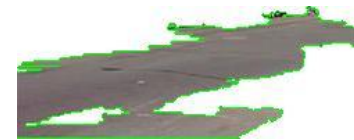
Step 2: Region-level matching



Pixel Area (size)



Road



Tree



Sky



Building



Snow

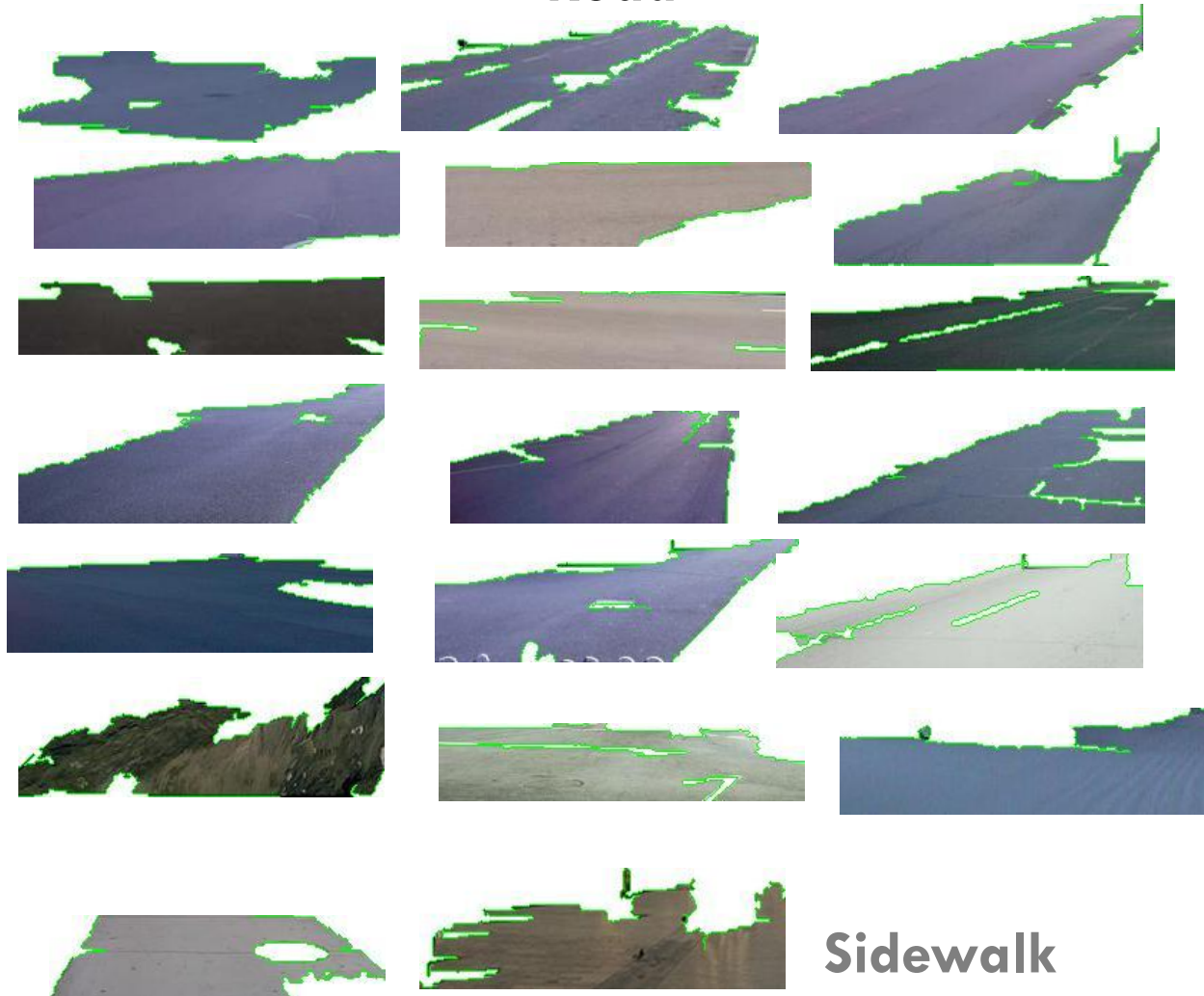
Step 2: Region-level matching



Absolute mask
(location)



Road

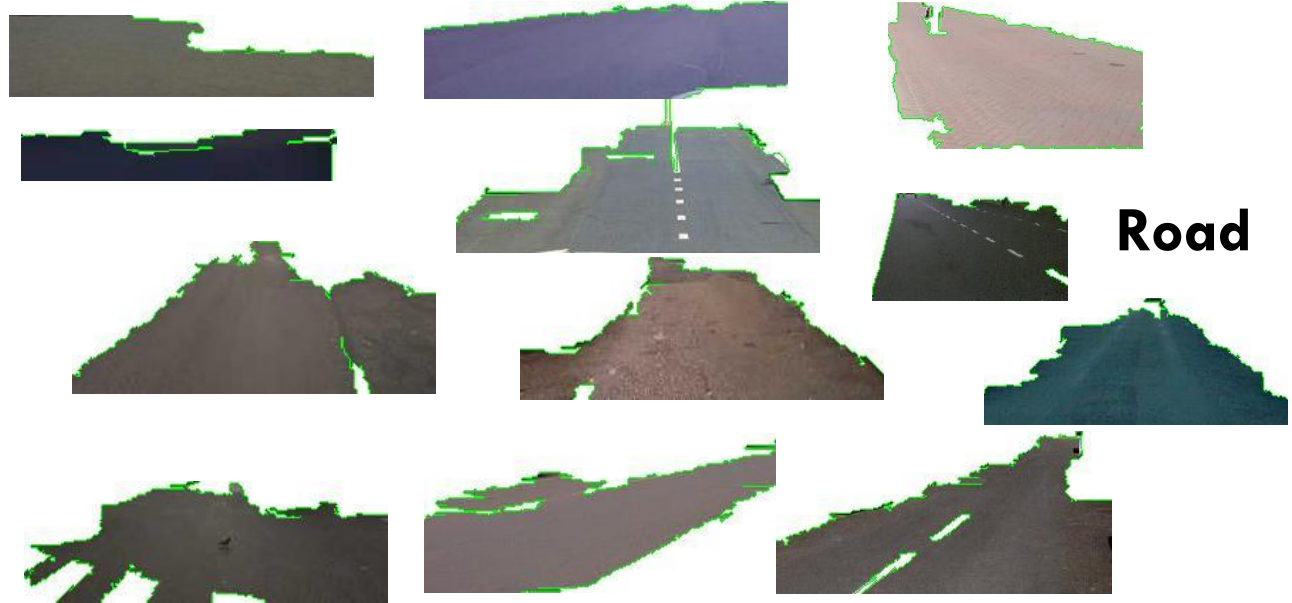


Sidewalk

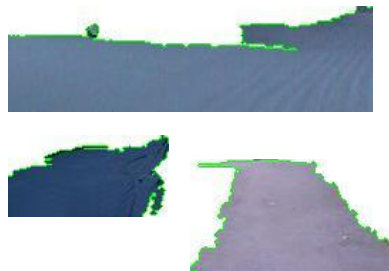
Step 2: Region-level matching



Texture



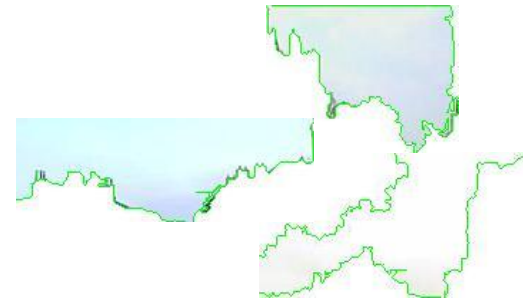
Road



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

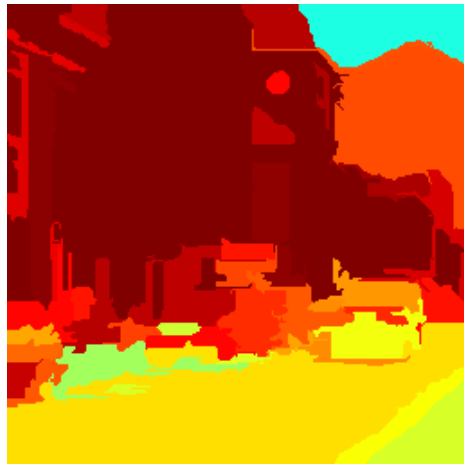
Features of class c within some radius of r_i

Total features of class c in the dataset

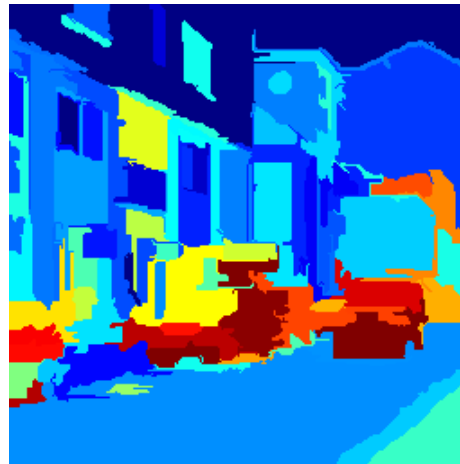
- Per-feature likelihoods combined via Naïve Bayes:

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

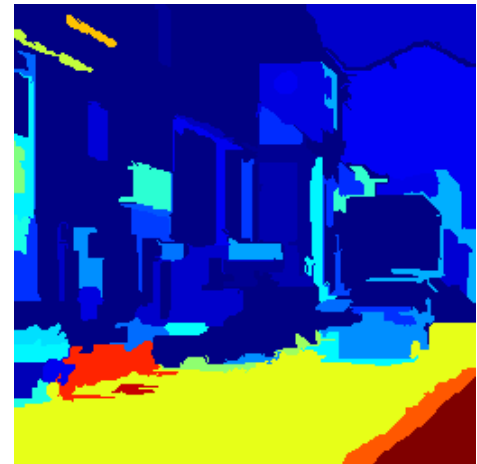
Region-level likelihoods



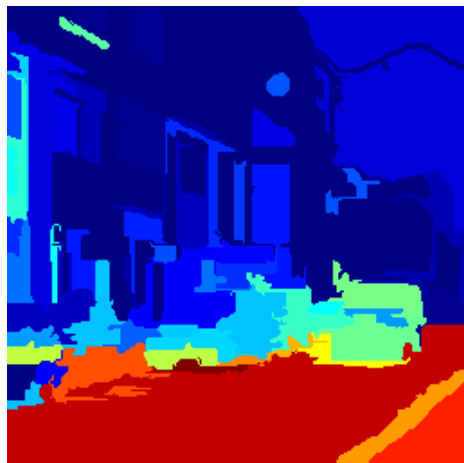
Building



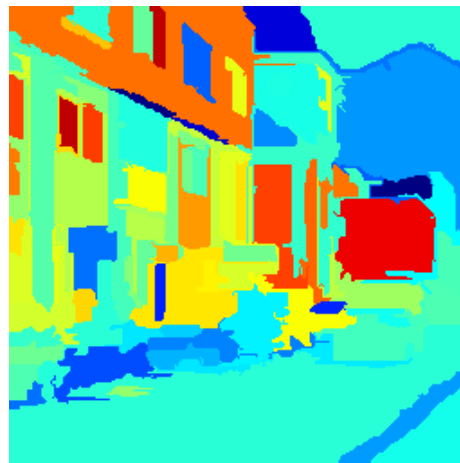
Car



Crosswalk



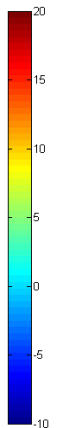
Road



Window



Sky



Step 3: Global image labeling

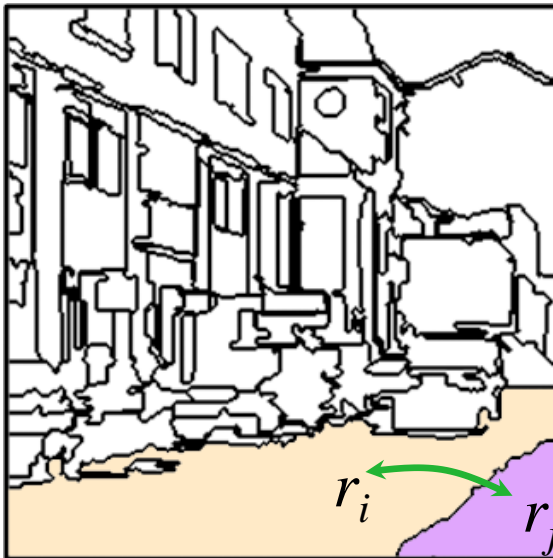
- Compute a global image labeling by optimizing a Markov random field (MRF) energy function:

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

↑
Vector of region labels

Regions

Neighboring regions



Efficient approximate minimization using α -expansion (Boykov et al., 2002)

Step 3: Global image labeling

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

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

Diagram annotations:

- Vector of region labels (points to \mathbf{c})
- Regions (points to i)
- Neighboring regions (points to i, j)

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

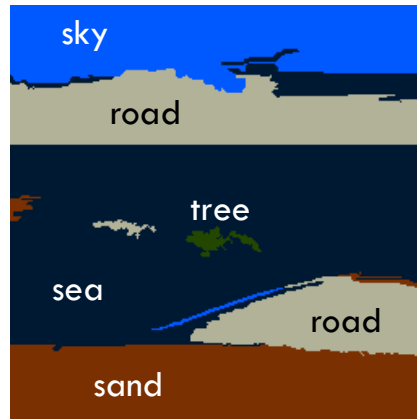
Regions

Neighboring regions

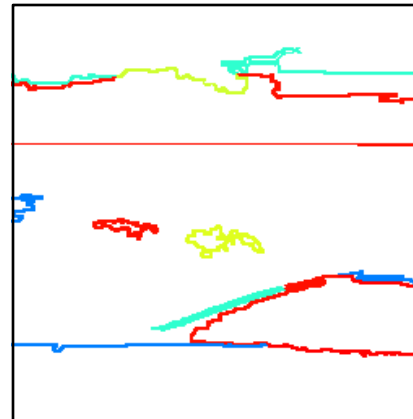
Original image



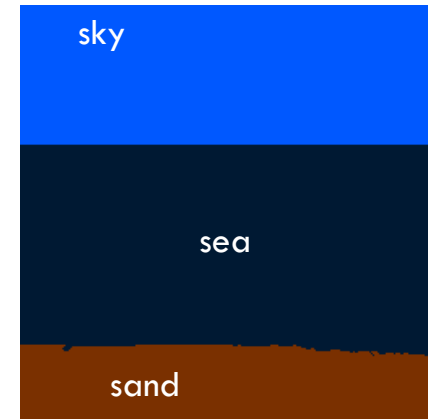
Maximum likelihood labeling



Edge penalties

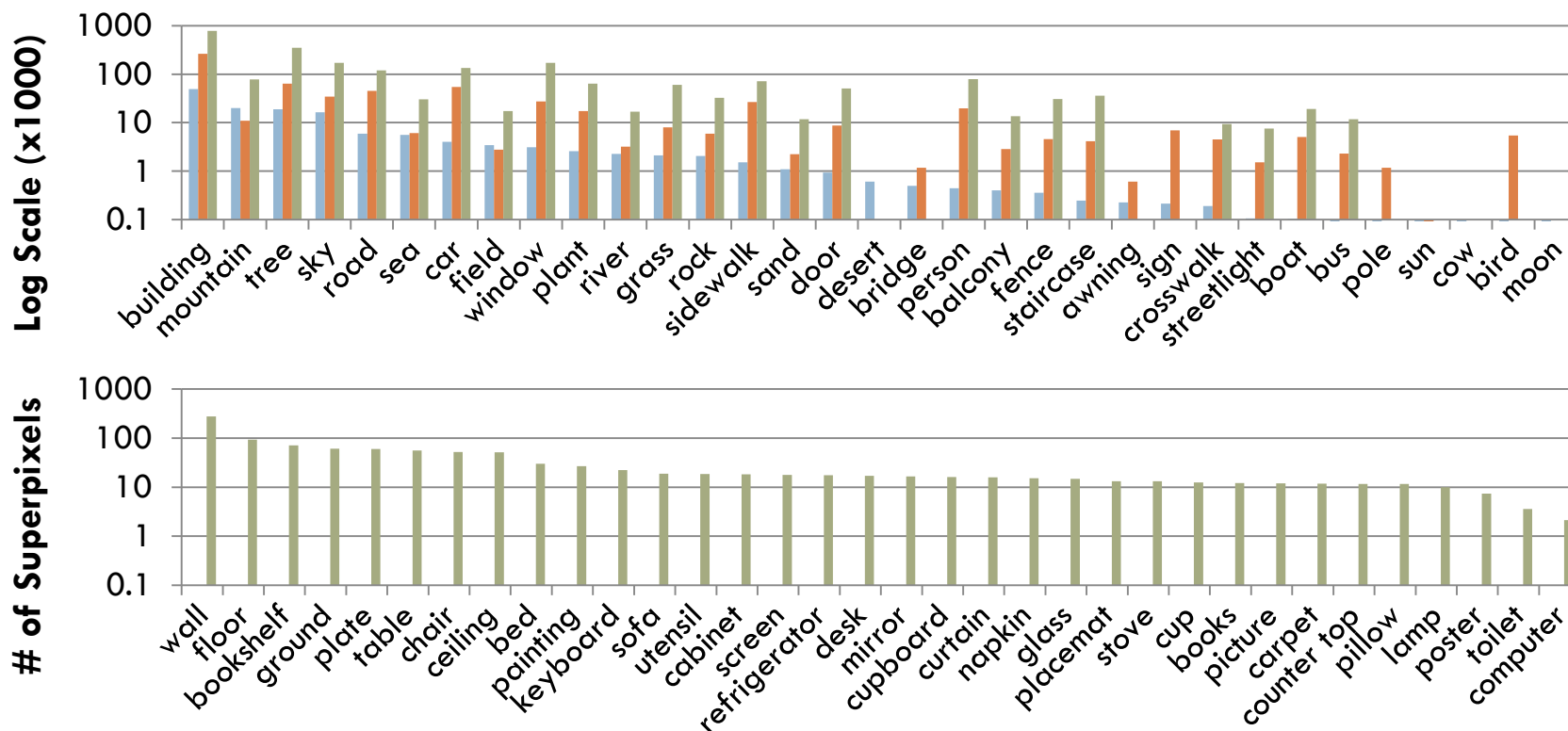


MRF labeling

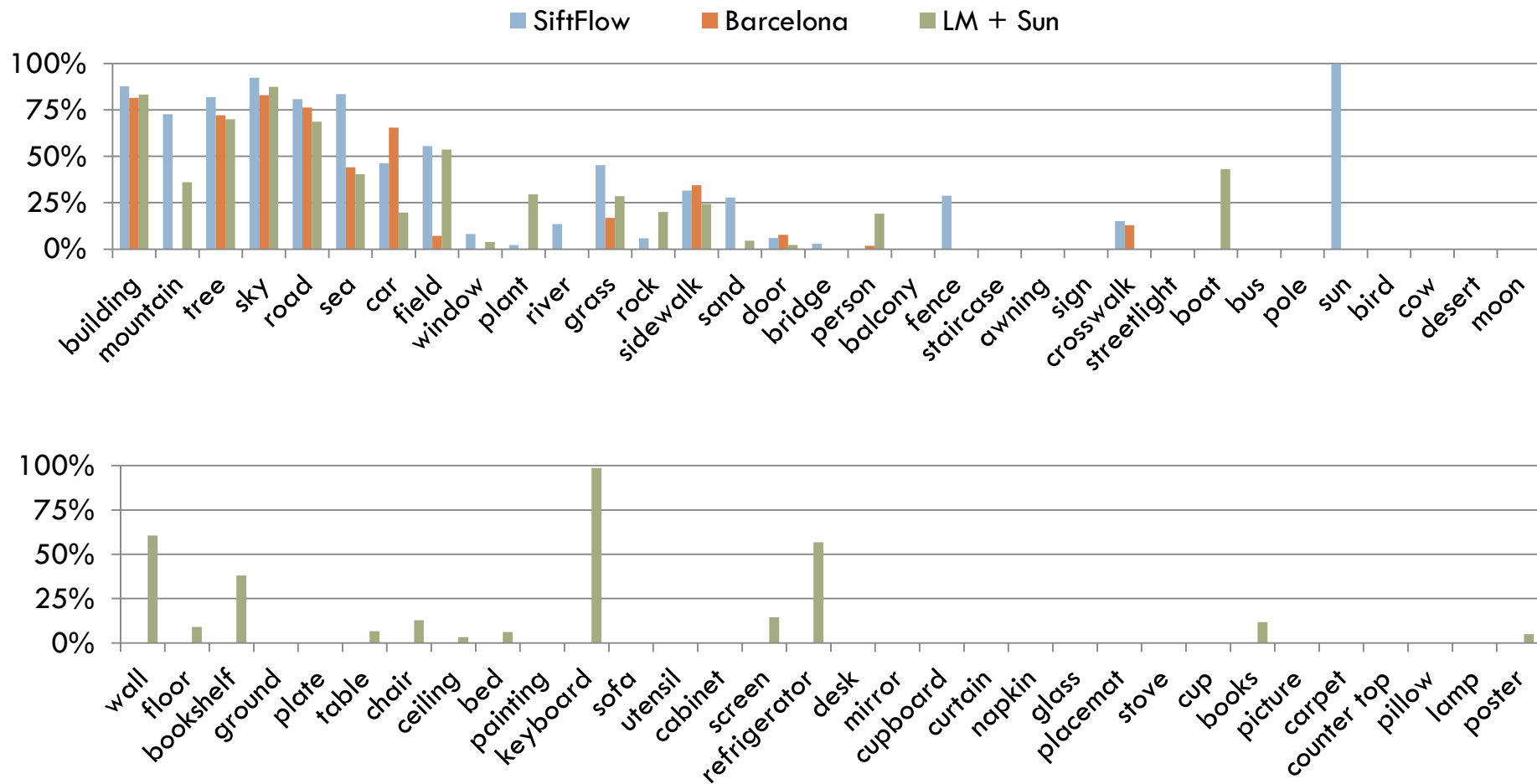


Datasets

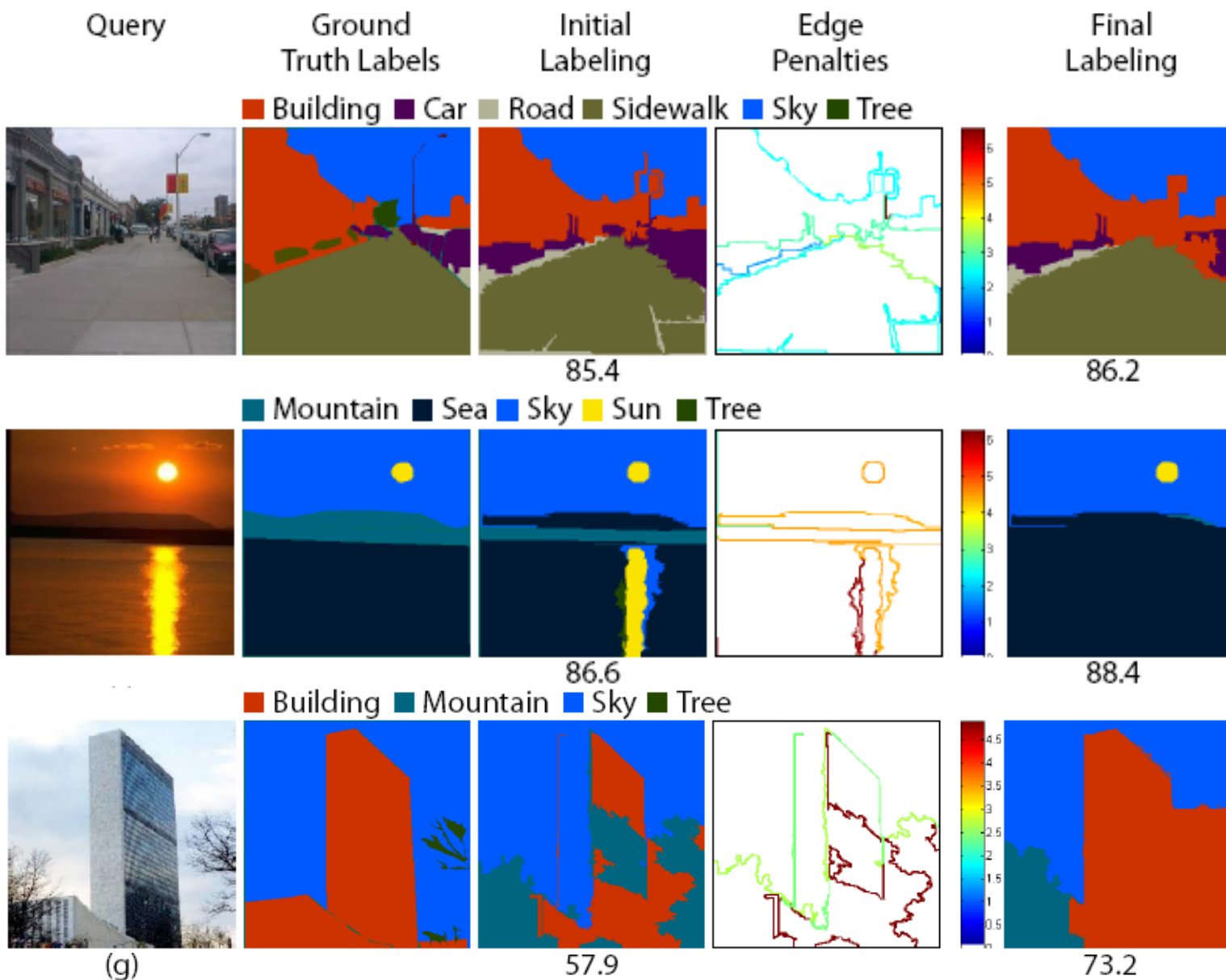
	Training images	Test images	Labels
SIFT Flow (Liu et al., 2009)	2,488	200	33
Barcelona	14,871	279	170
LabelMe+SUN	50,424	300	232



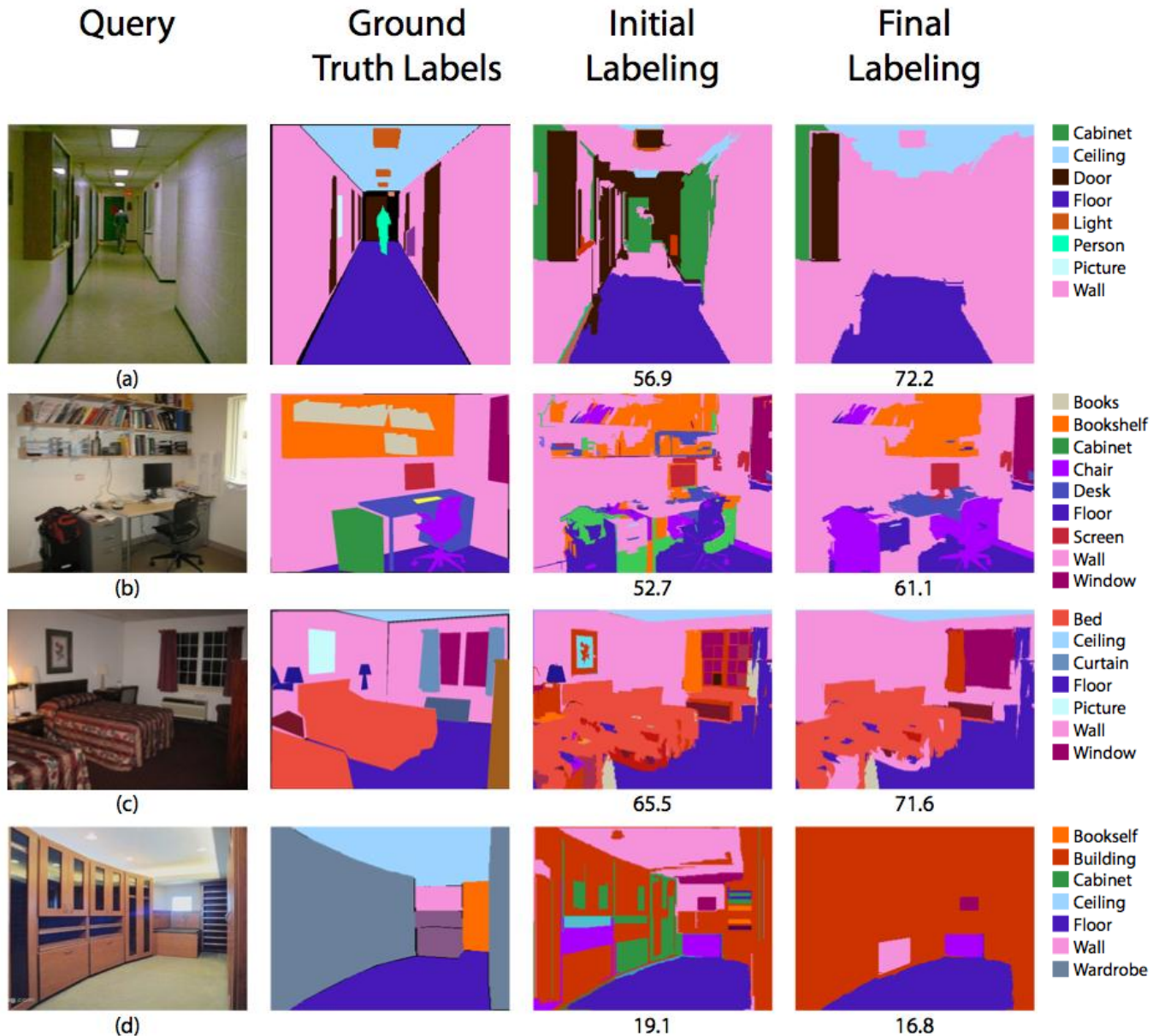
Per-class classification rates



Results on SIFT Flow dataset



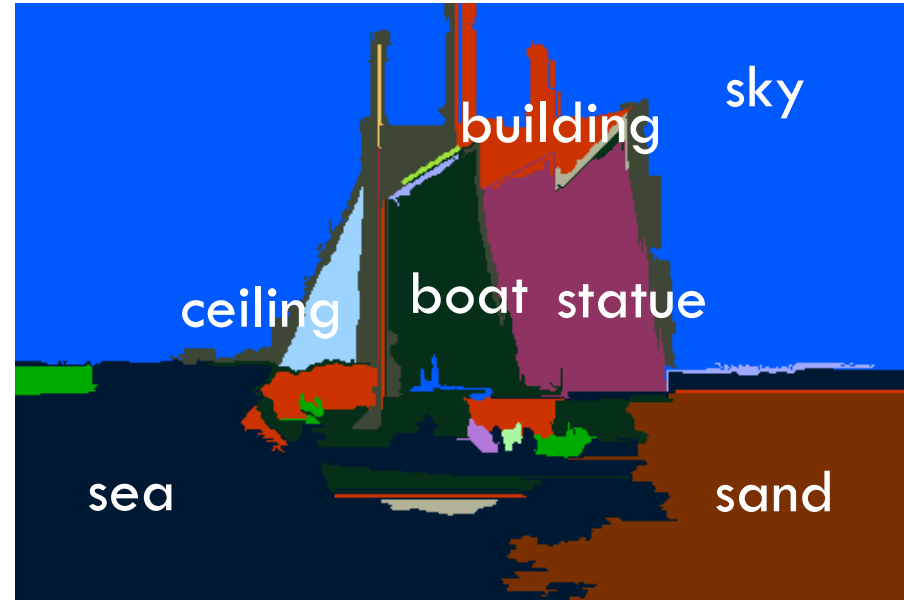
Results on LM+SUN dataset



Summary so far

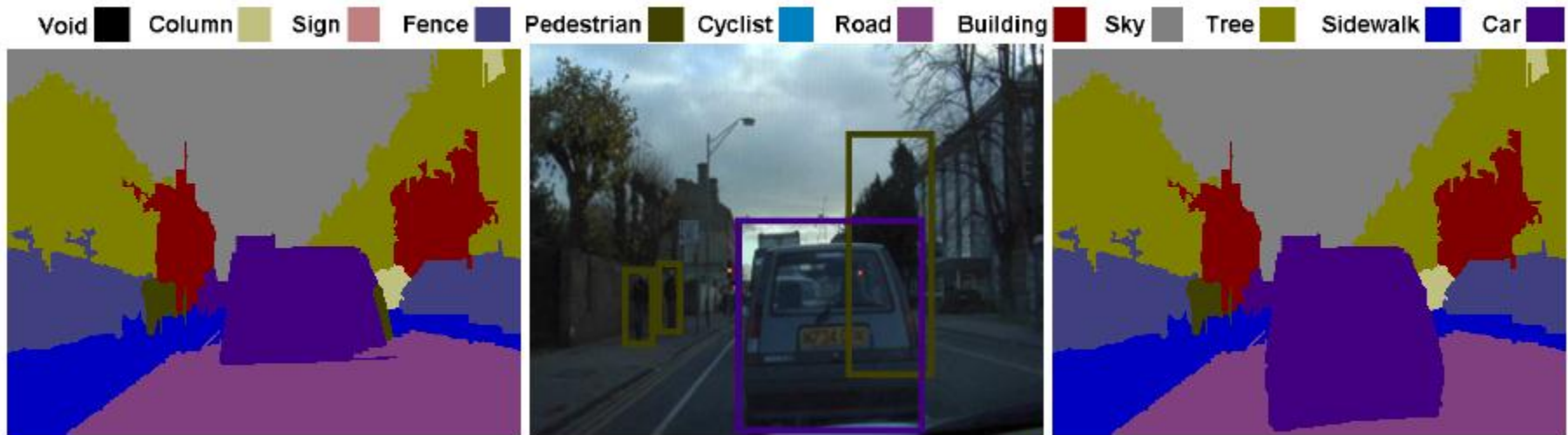
- A lazy learning method for image parsing:
 - ▣ Global scene matching
 - ▣ Superpixel-level matching
 - ▣ MRF optimization
- Challenges
 - ▣ Indoor images are hard!
 - ▣ We do well on “stuff” but not on “things”

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



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

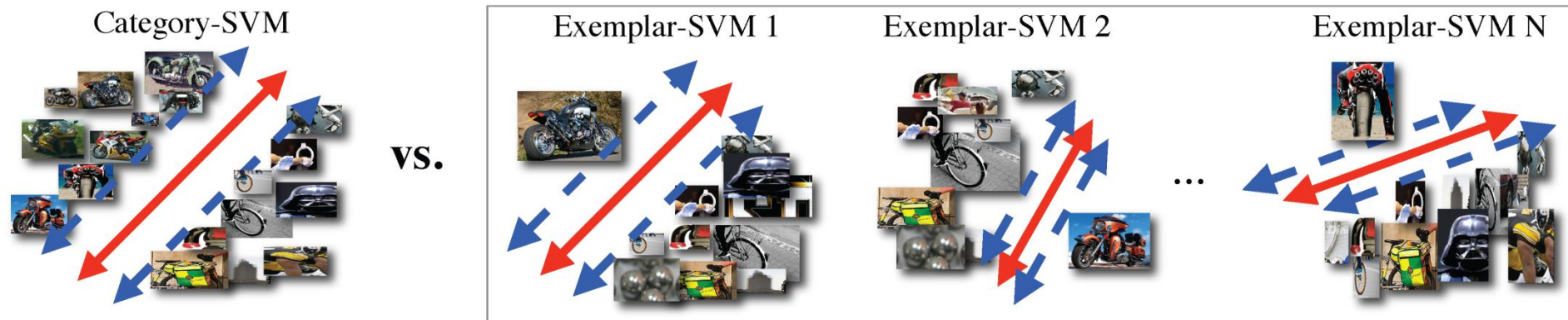
Problems with this approach

- 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

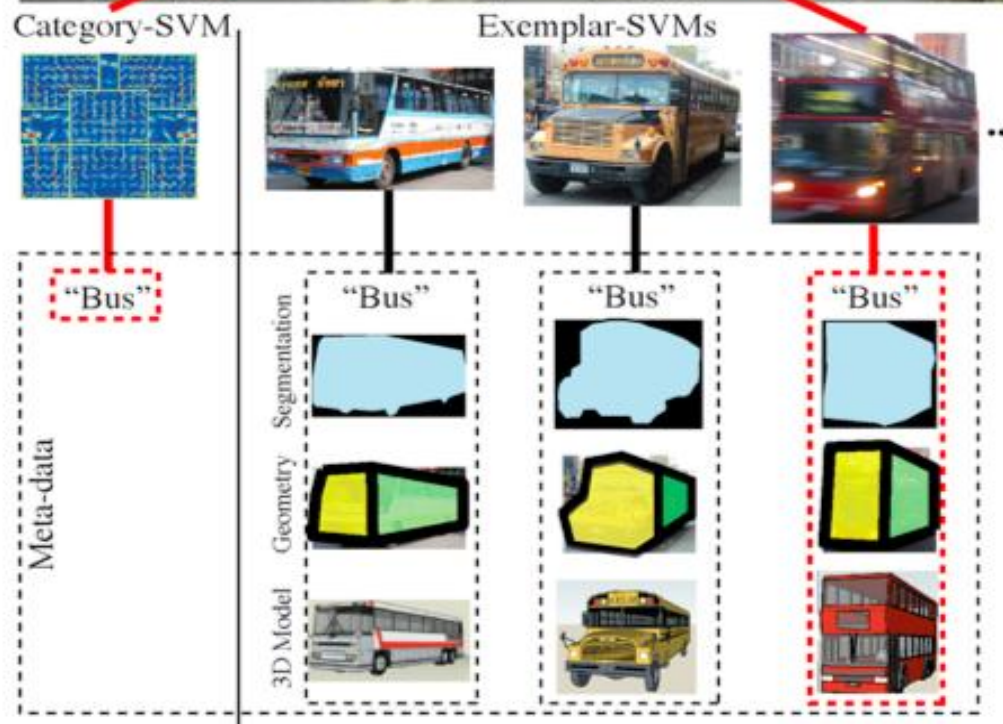
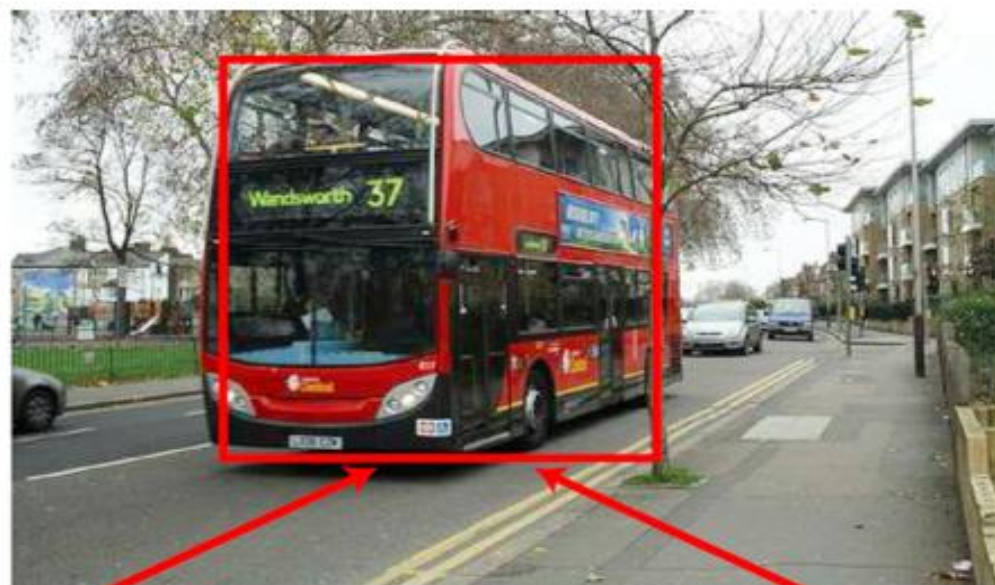


Per-exemplar detectors

- 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



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

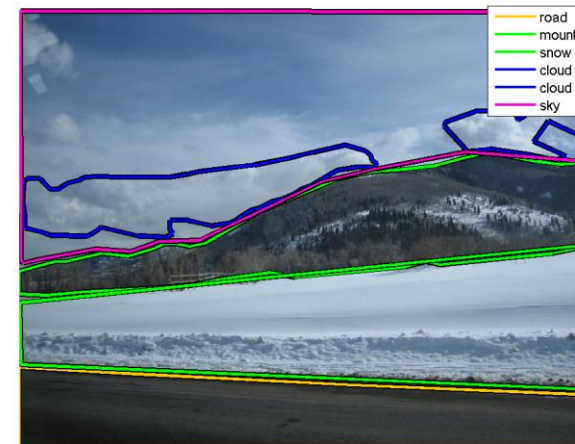
Retrieval set for



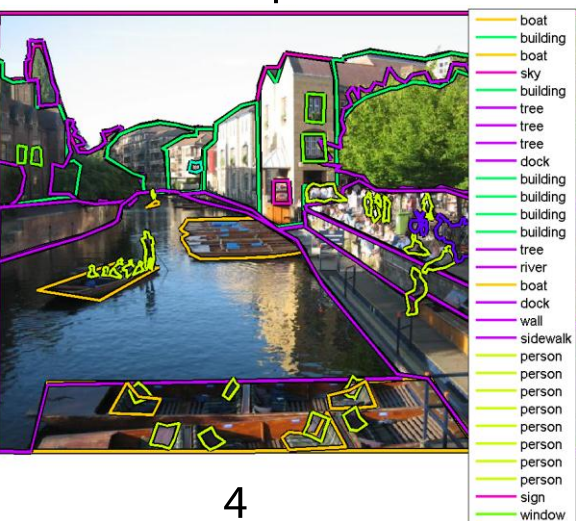
1



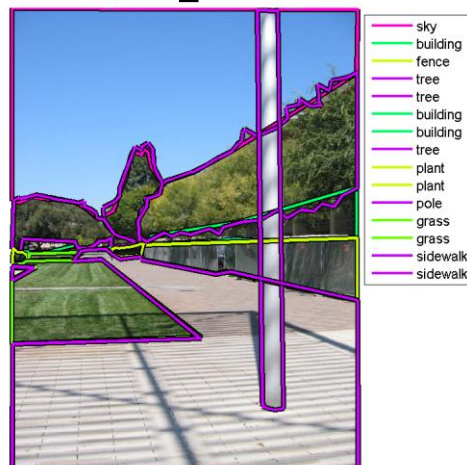
2



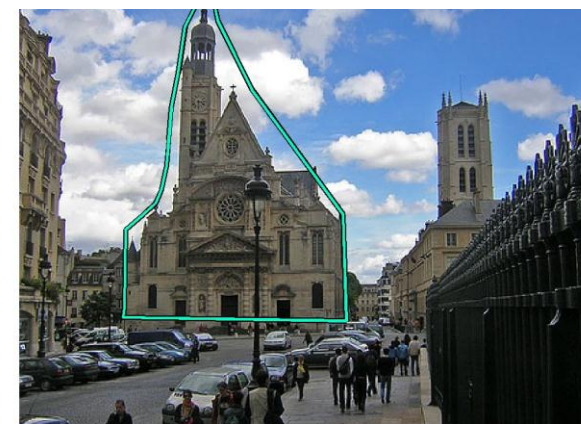
3



4

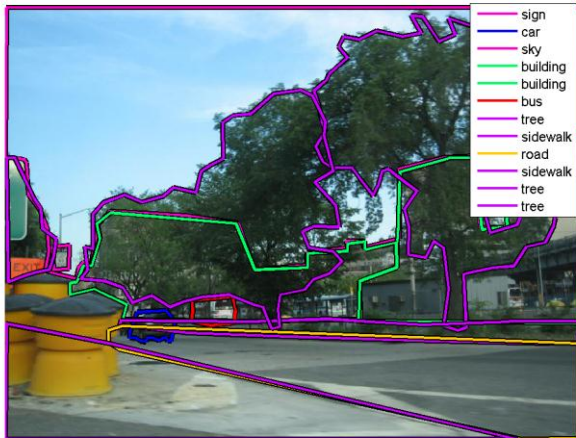


5

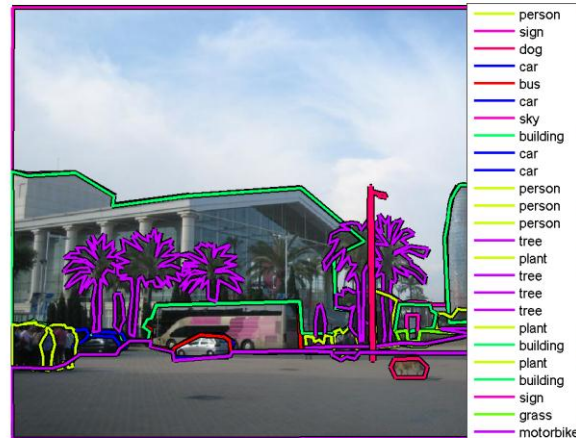


6

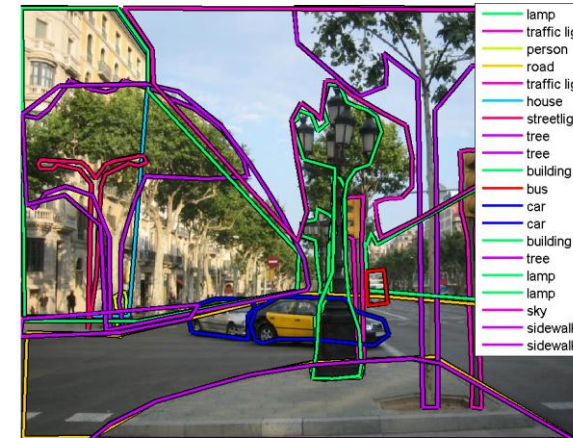
Retrieval set for



16



26



59



342



410



491

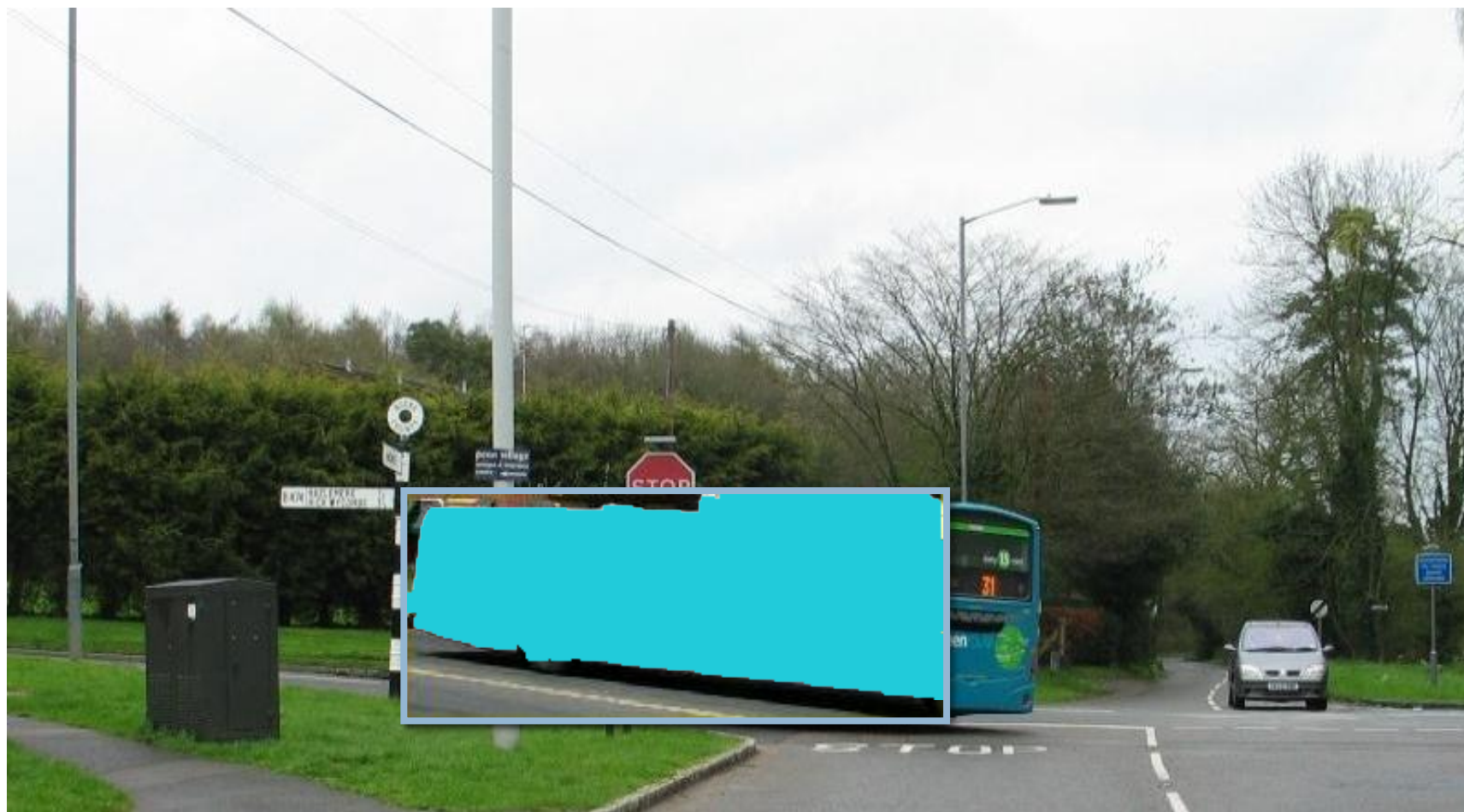
Per-exemplar detectors for parsing

- 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 masked for all positive detections

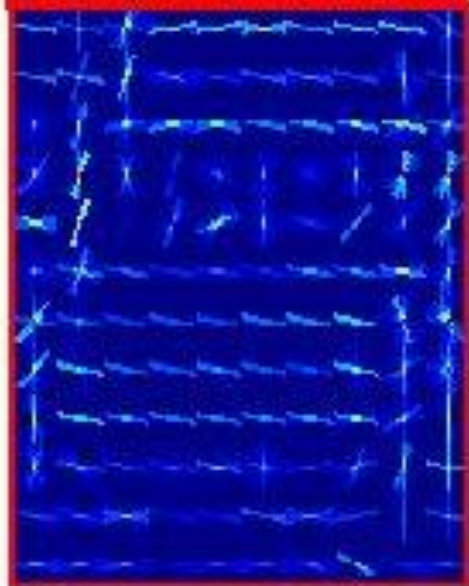
Per-exemplar detectors for parsing



Per-exemplar detectors for parsing



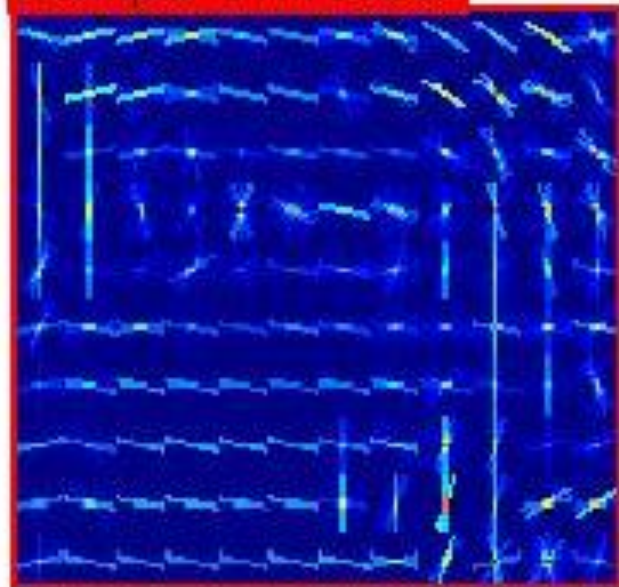
Exemplar-SVM bus 20



Exemplar Image 20



Exemplar-SVM bus 28



Exemplar Image 28

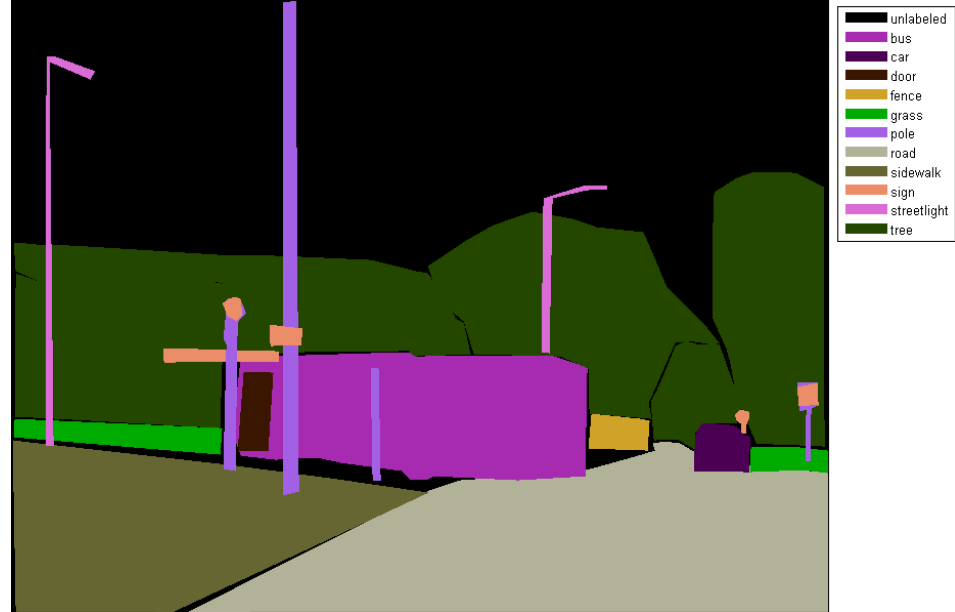


Per-exemplar detectors for parsing

- 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



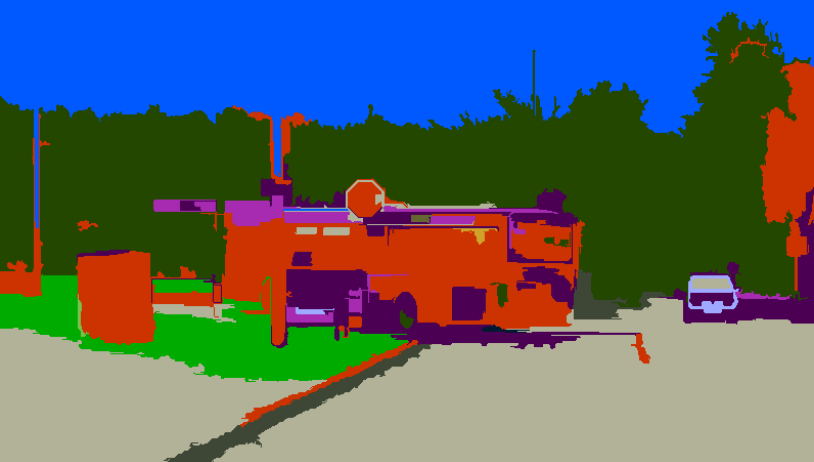
55% (23%)

Detector-based Parsing Result



45% (26%)

Superparsing Result



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

55% (23%)



Detector Based Parsing Result



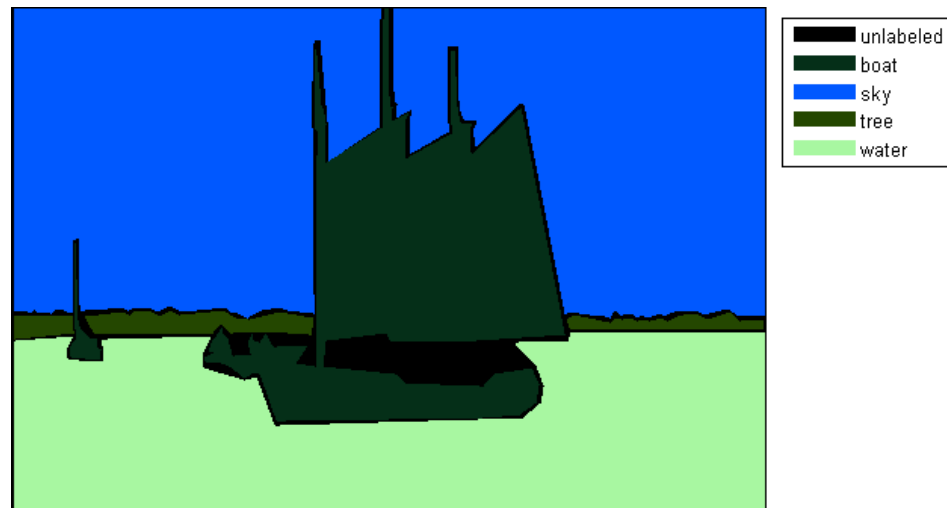
- air conditioner
- awning
- basket
- bench
- boat
- bowl
- box
- building
- bus
- bush
- cabinet
- car
- column
- counter top
- cup
- cupboard
- door
- fence
- fish
- fountain
- glass
- grass
- ground
- handrail
- hill
- house
- jar
- laptop
- leg
- mountain
- mousepad

45% (26%)



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

61% (31%)



Superparsing Result



52% (31%)

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

Detector Based Parsing Result



19% (25%)

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

Superparsing Result



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

52% (31%)

Detector Based Parsing Result



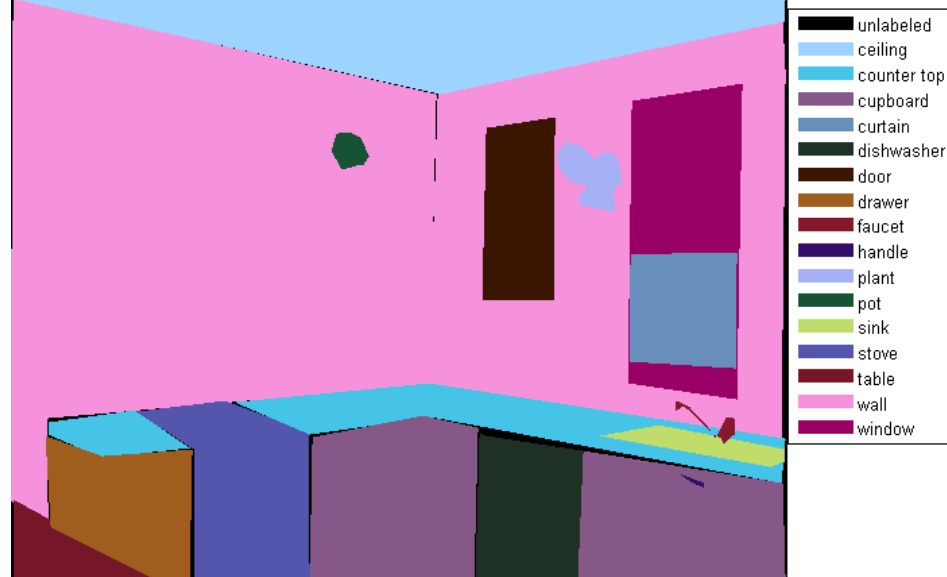
- air conditioner
- airplane
- boat
- books
- bookshelf
- bridge
- building
- car
- ceiling
- door
- field
- grass
- ground
- hill
- mountain
- pen
- plate

19% (25%)

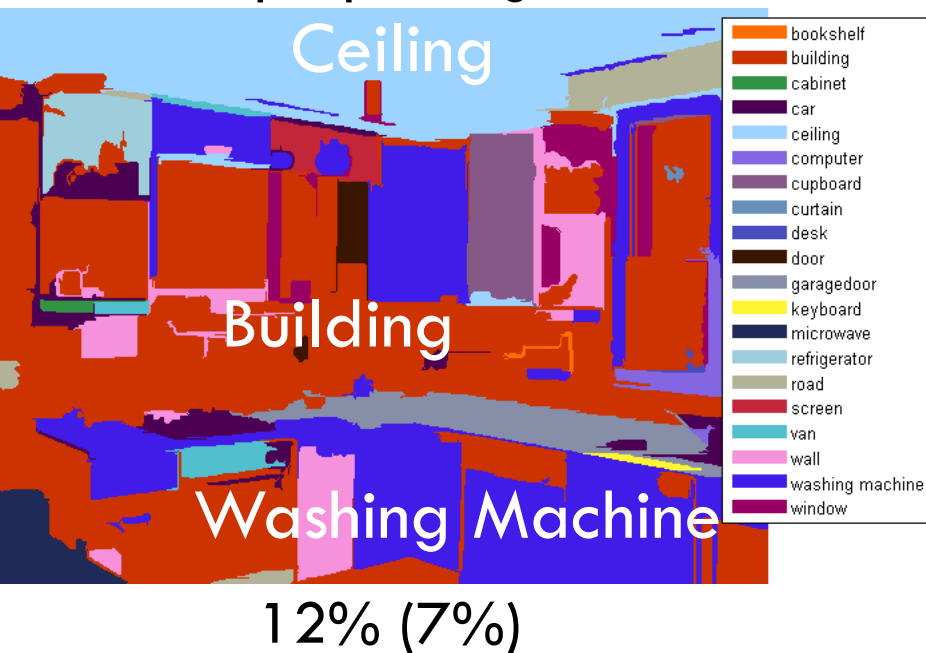


- boat
- building
- church
- grass
- mountain
- road
- sand
- sea
- sky
- wall

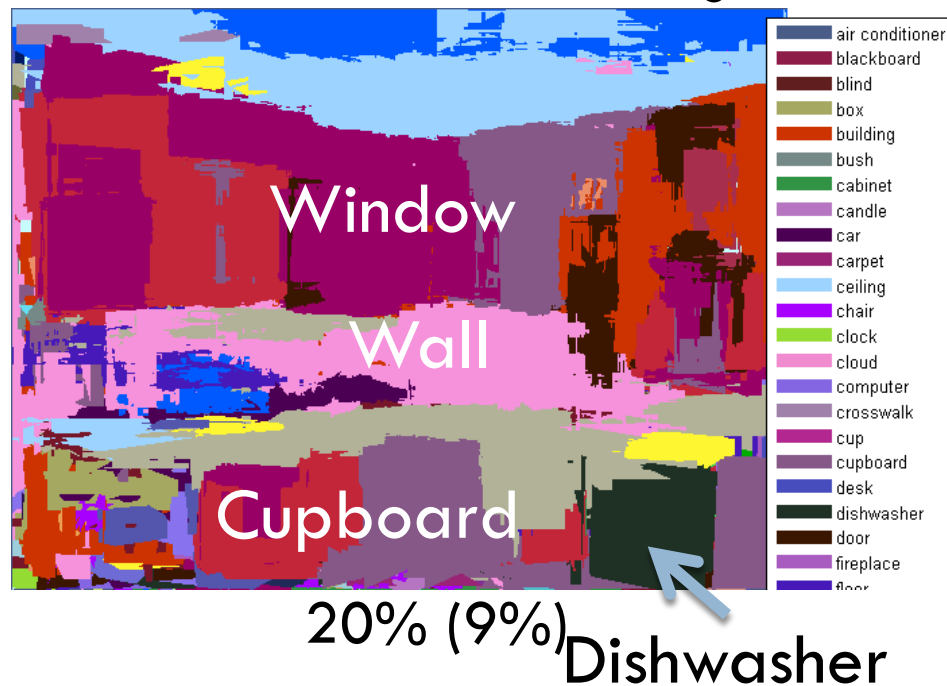
62% (46%)



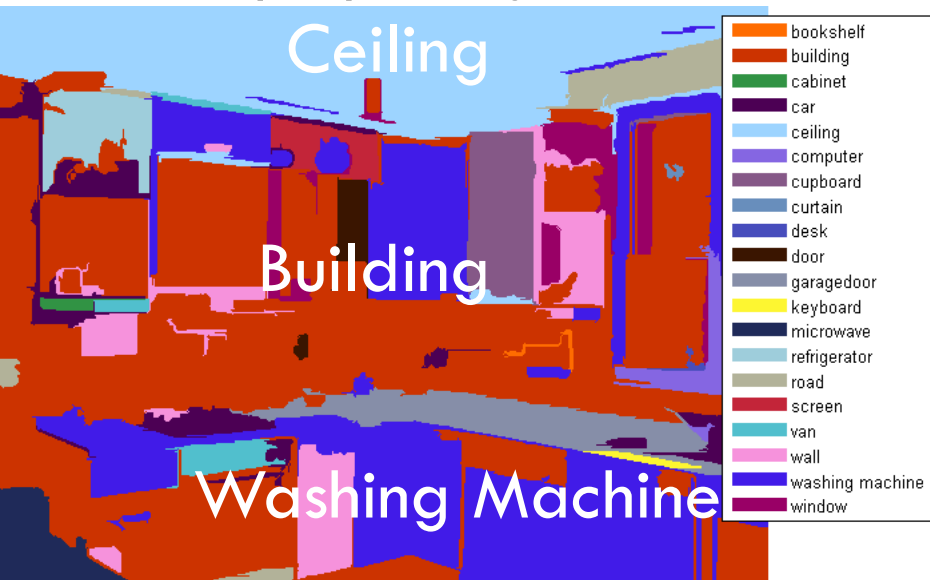
Superparsing Result



Detector Based Parsing Result



Superparsing Result



12% (7%)

Detector Based Parsing Result



20% (9%)



24% (10%)

Next Steps

- Determine which detectors to run
 - ▣ Manually select the “thing” classes
 - ▣ Use context from parsing
- Find better methods for integrating image parsing and detectors

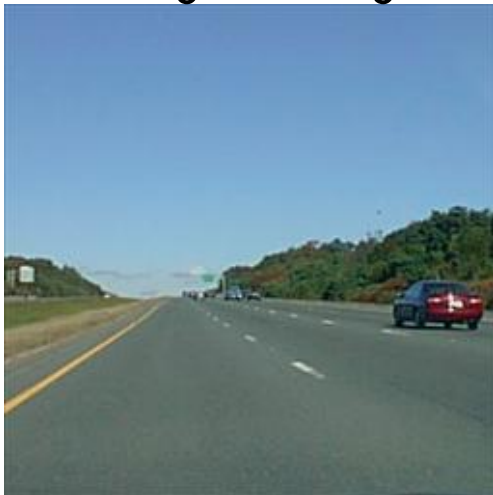
Review so far

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

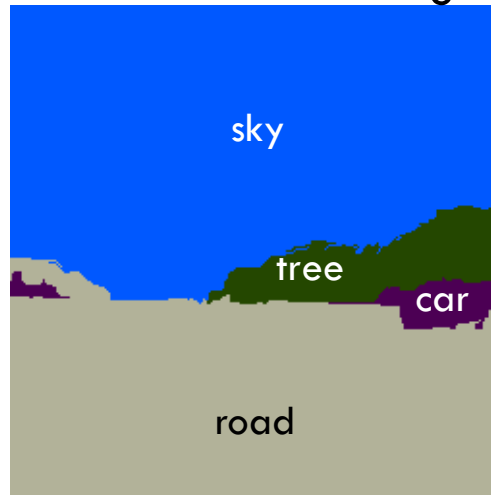
Joint geometric/semantic labeling

- **Semantic labels:** road, grass, building, car, etc.
- **Geometric labels:** sky, vertical, horizontal
 - ▣ Gould et al. (ICCV 2009)

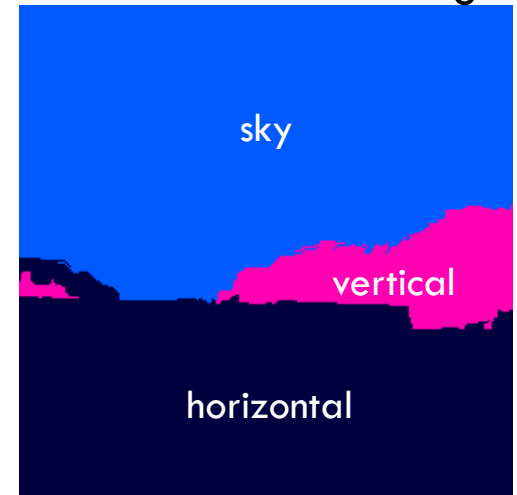
Original image



Semantic labeling



Geometric labeling



Recall: Global image labeling

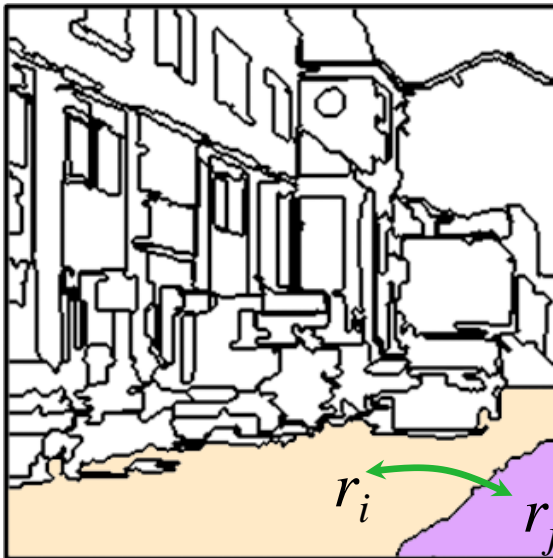
- Compute a global image labeling by optimizing a Markov random field (MRF) energy function:

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

↑
Vector of region labels

Regions

Neighboring regions



Efficient approximate minimization using α -expansion (Boykov et al., 2002)

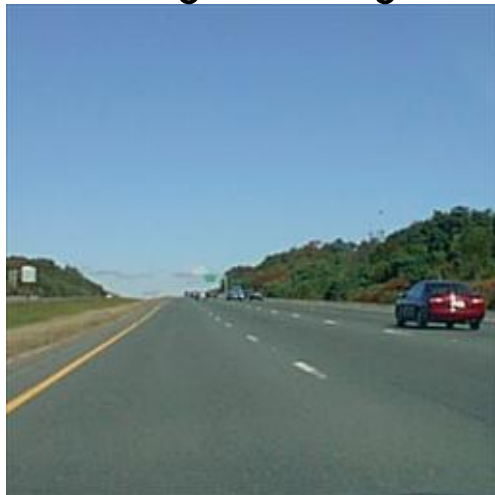
Joint geometric/semantic labeling

- Objective function for joint labeling:

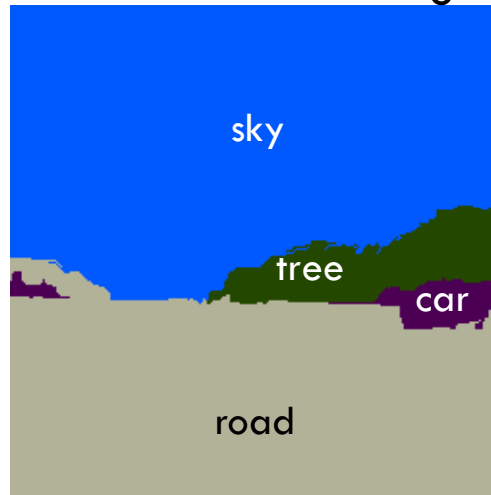
$$F(\mathbf{c}, \mathbf{g}) = \underbrace{E(\mathbf{c})}_{\text{Cost of semantic labeling}} + \underbrace{E(\mathbf{g})}_{\text{Cost of geometric labeling}} + \mu \sum_{\text{regions } r_i} \underbrace{\psi(c_i, g_i)}_{\text{Geometric/semantic consistency penalty}}$$

Semantic labels *Geometric labels*

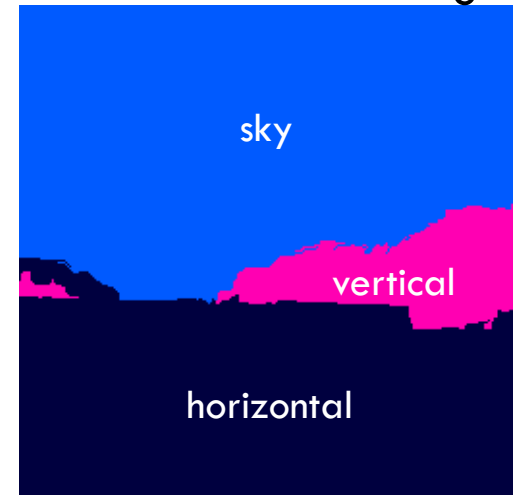
Original image



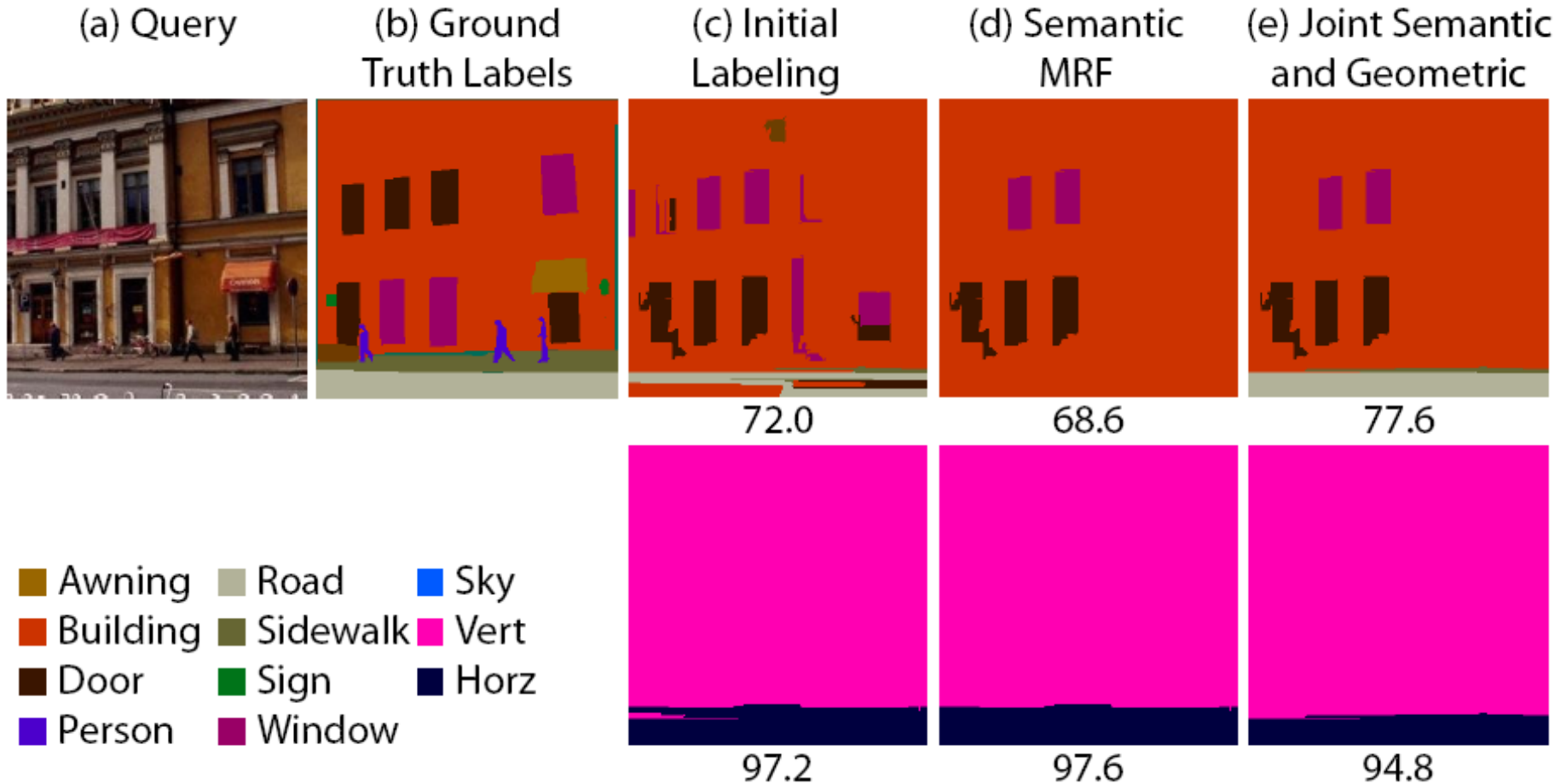
Semantic labeling



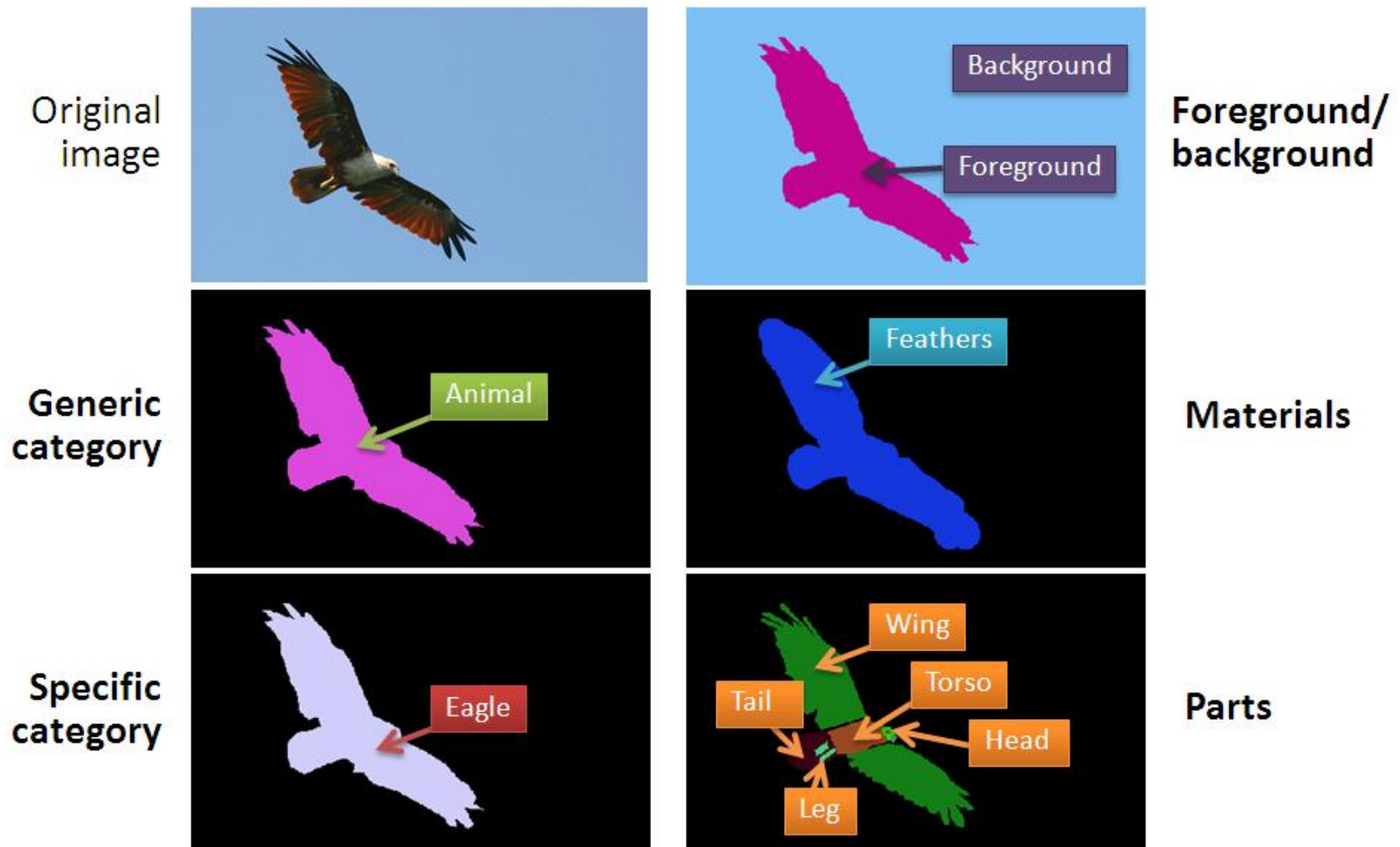
Geometric labeling



Example of joint labeling



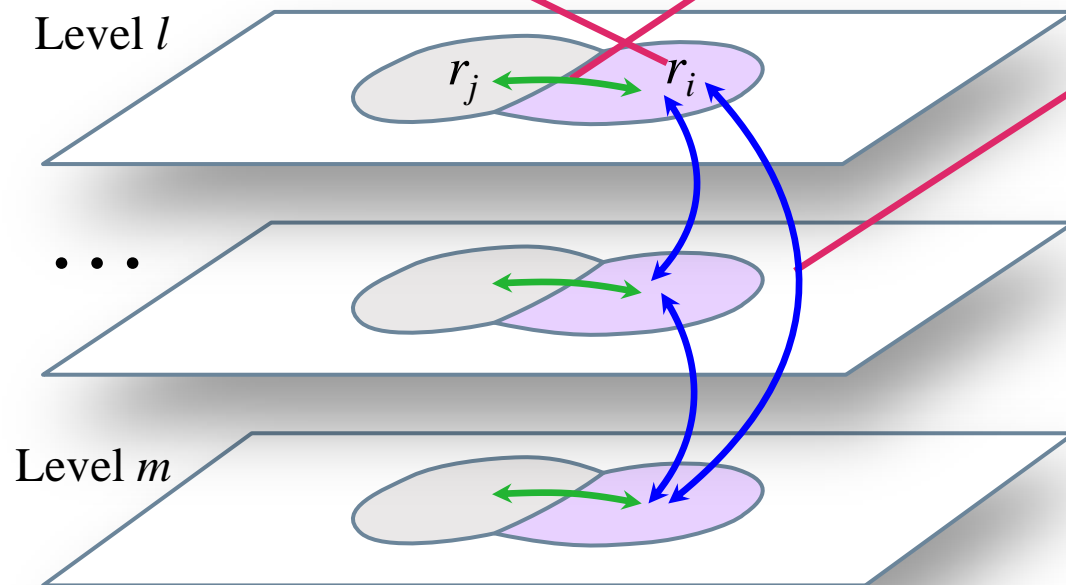
Understanding scenes on many levels



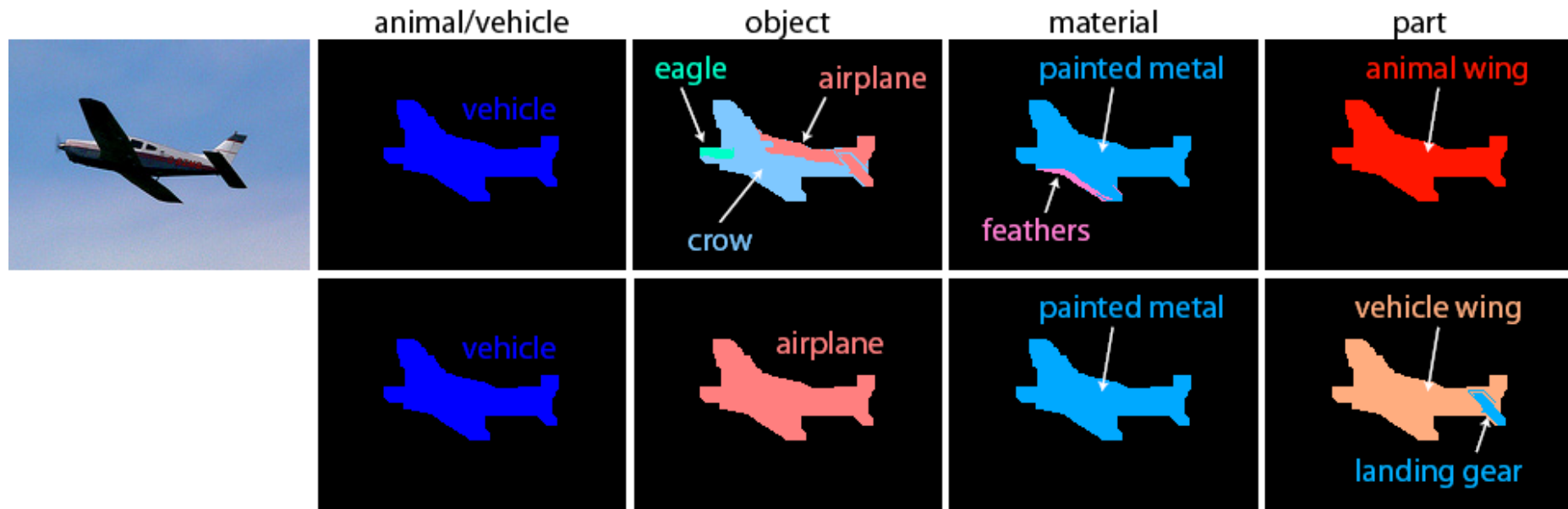
Understanding scenes on many levels

$$E(\mathbf{c}^1, \dots, \mathbf{c}^L) = \sum_l \left(\underbrace{\sum_i E_{\text{data}}^l(r_i, c_i^l)}_{\text{Likelihood score for level } l, \text{ region } r_i \text{ and label } c_i} + \underbrace{\sum_{i,j} E_{\text{horiz}}^l(c_i^l, c_j^l)}_{\text{Single-level MRF Consistency penalty for neighbors on the same level}} \right) + \sum_{l \neq m} \sum_i \underbrace{E_{\text{vert}}^{l,m}(c_i^l, c_i^m)}_{\text{Consistency penalty for labels of region } i \text{ on levels } l \text{ and } m}$$

Region labelings for every level



Understanding scenes on many levels



Review



- Nonparametric image parsing
- Beyond superpixels
- Beyond unique labels

Review so far

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

Beyond labels: Attributes

