



# Multiple View Semantic Segmentation for Street View Images

Jianxiong Xiao and Long Quan

{csxjx, quan}@cse.ust.hk

The Hong Kong University of Science and Technology

## Abstract

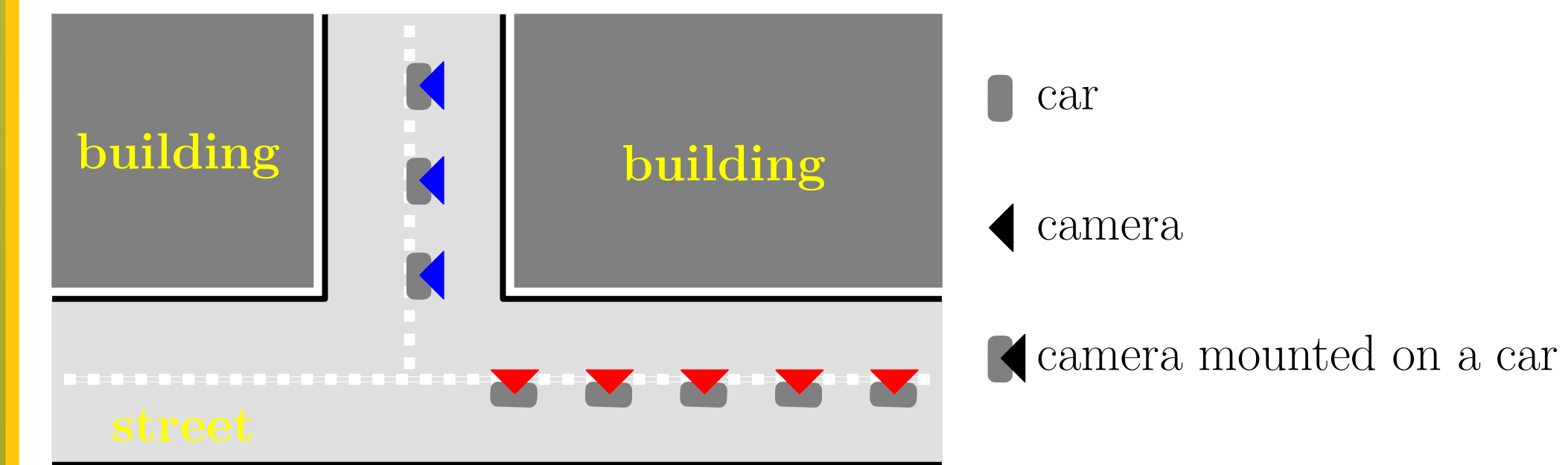
We propose a simple but powerful multi-view semantic segmentation framework for images captured by a camera mounted on a car driving along streets.

In our approach, a pair-wise Markov Random Field (MRF) is laid out across multiple views. Both 2D and 3D features are extracted at a super-pixel level to train classifiers for the unary data terms of MRF. For smoothness terms, our approach makes use of color differences in the same image to identify accurate segmentation boundaries, and dense pixel-to-pixel correspondences to enforce consistency across different views. To speed up training and to improve the recognition quality, our approach adaptively selects the most similar training data for each scene from the label pool.

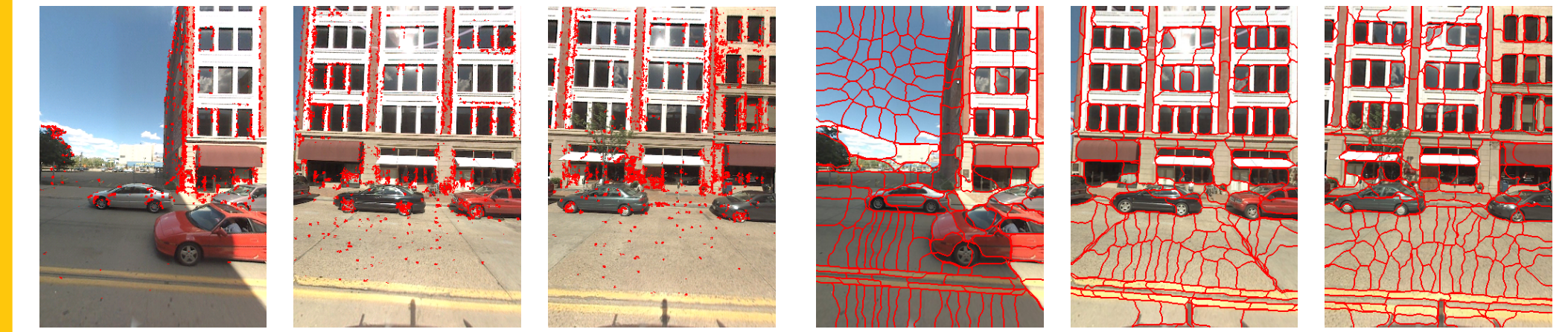
Furthermore, we also propose a powerful approach within the same framework to enable large-scale labeling in both the 3D space and 2D images. We demonstrate our approach on more than 10,000 images from Google Maps Street View.

## Preprocessing

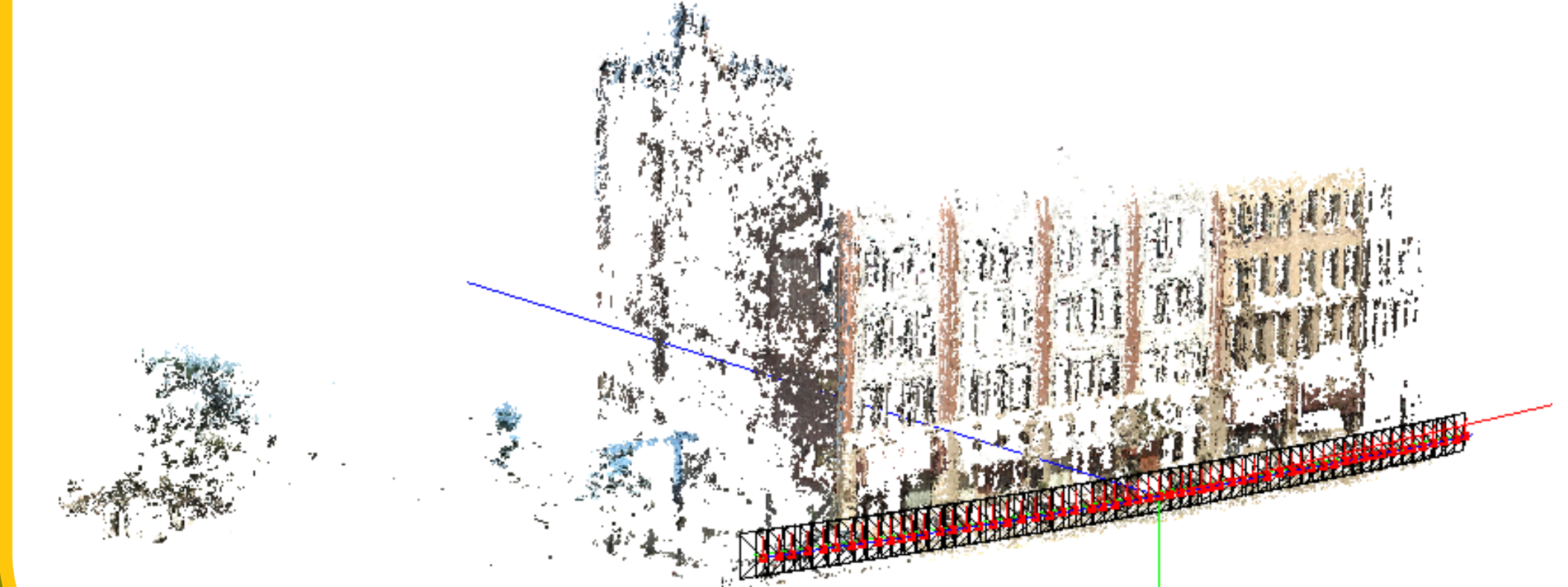
Data is captured by a camera that usually faces the building façade and moves laterally along streets.



We compute pixel-to-pixel correspondences to obtain feature tracks across multiple view, and over-segment the images into superpixels.

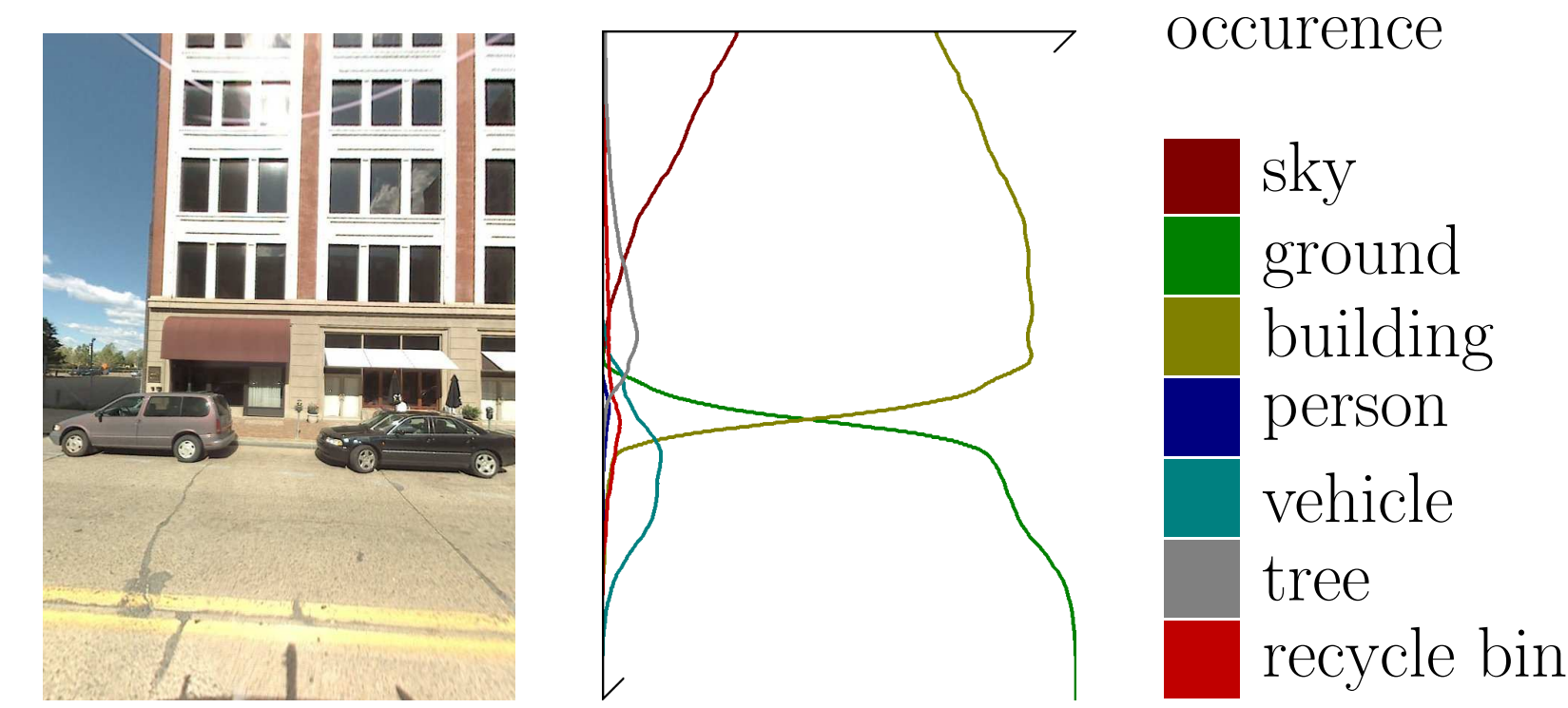


We use Structure from Motion to obtain a 3D point cloud for the scene.

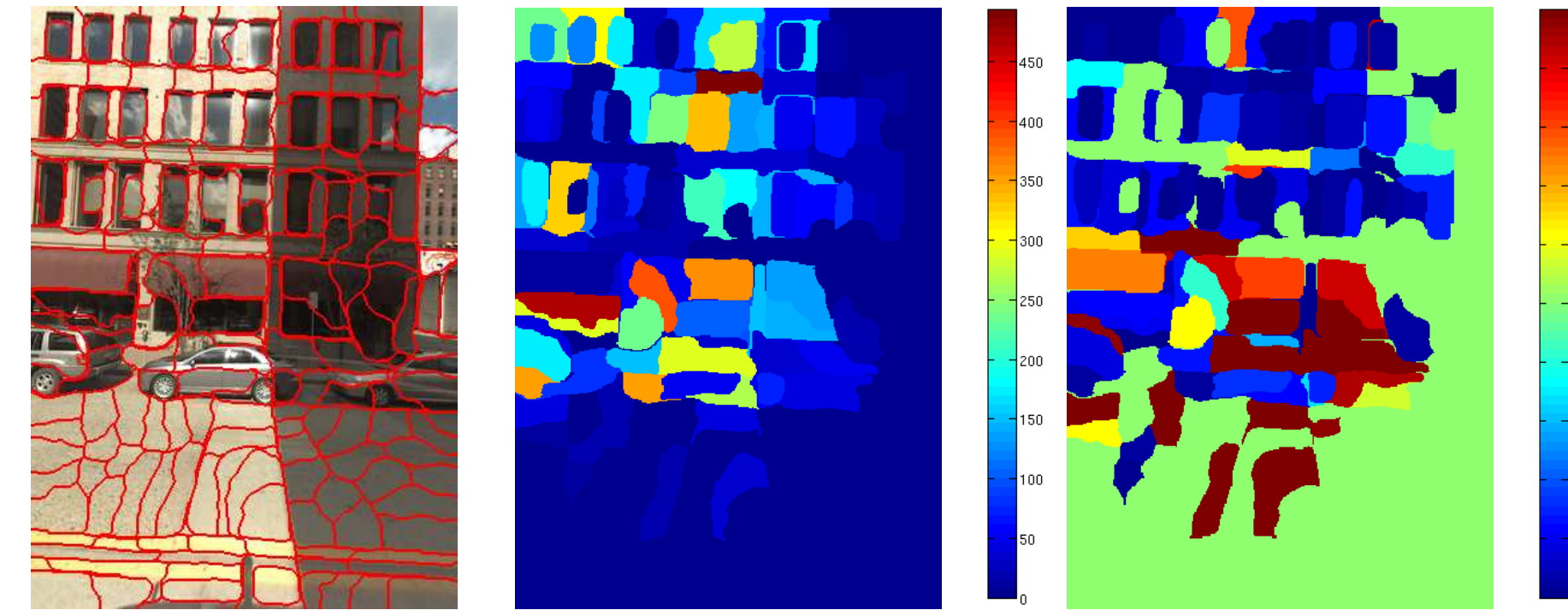


## Unary Potential

For each superpixel, we extract 2D features from texture, color, size, shape and pixel position in the image.



We extract 3D features as the feature track density in SFM and the dot product between the estimated normal direction and  $-y$  direction.



We apply the AdaBoost classifier that we have learned for each class,  $l$ , to the descriptors. The estimated confidence value can be reinterpreted as a probability distribution using softmax transformation:

$$P_i(l_i | \mathbf{f}_i^A, \mathbf{f}_i^P, \mathbf{f}_i^G) = \frac{\exp(H(l_i, \mathbf{f}_i^A, \mathbf{f}_i^P, \mathbf{f}_i^G))}{\sum_l \exp(H(l, \mathbf{f}_i^A, \mathbf{f}_i^P, \mathbf{f}_i^G))},$$

where  $H(l, \mathbf{f}_i^A, \mathbf{f}_i^P, \mathbf{f}_i^G)$  is the output of the AdaBoost classifier for class  $l$ . We then define the unary potential as  $\psi_i(l_i) = -\log P_i(l_i | \mathbf{f}_i^A, \mathbf{f}_i^P, \mathbf{f}_i^G)$ .

## Smoothness

For edges in the same image, the smoothness cost is defined as

$$\psi_{ij}(l_i, l_j) = [l_i \neq l_j] \cdot g(i, j),$$

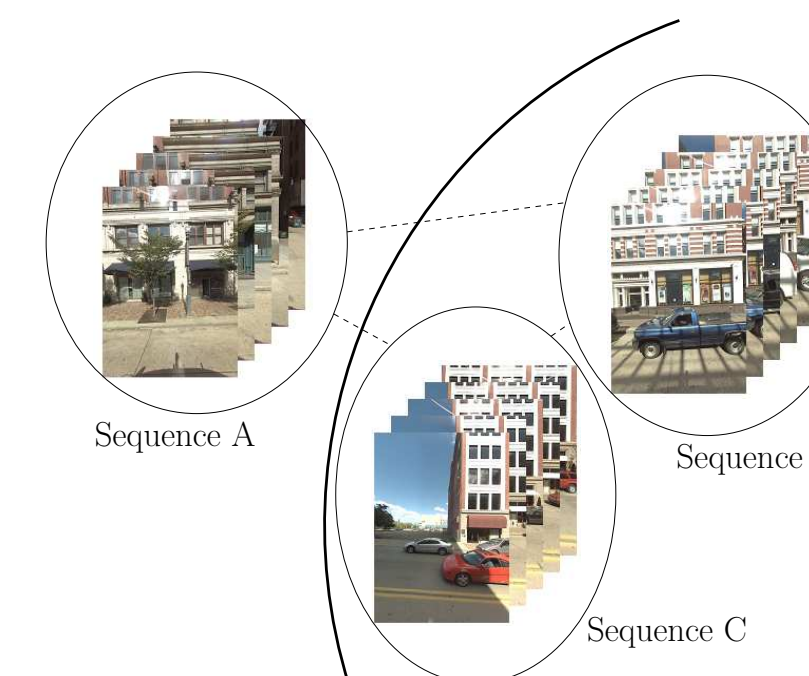
where  $g(i, j) = \frac{1}{\zeta \|\mathbf{c}_i - \mathbf{c}_j\|^2 + 1}$  and  $\|\mathbf{c}_i - \mathbf{c}_j\|^2$  is the  $L2$ -Norm of the RGB color difference of two superpixels,  $p_i$  and  $p_j$ .

For edge across two images, the smoothness cost is defined as

$$\psi_{ij}(l_i, l_j) = [l_i \neq l_j] \cdot \lambda |\mathcal{T}_{ij}| g(i, j)$$

where  $\mathcal{T}_{ij} = \{\mathbf{t} = \langle \mathbf{x}, (x_i, y_i, i), (x_j, y_j, j), \dots \rangle\}$  is the set of all feature tracks with projection  $(x_i, y_i)$  lying inside the superpixel,  $p_i$ , in image  $I_i$ , and projection  $(x_j, y_j)$  lying inside the superpixel,  $p_j$ , in image  $I_j$ . This definition encourages two superpixels with more feature track connections to have the same semantic segmentation label, since the cost to have different labels is high due to large  $|\mathcal{T}_{ij}|$ .

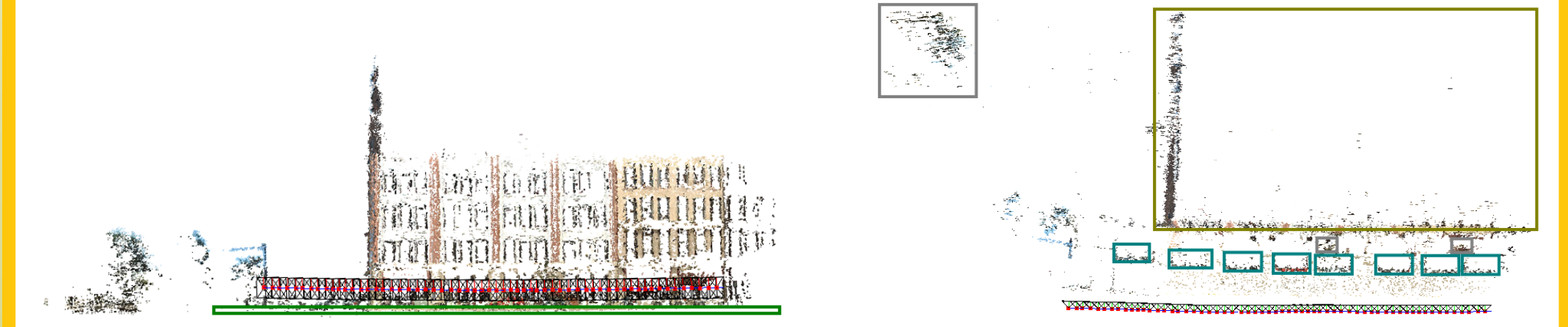
## Adaptive Training



We cluster the 40 labeled sequences in the pool based on affinity, defined as the minimal Gist distance between any image in a sequence and any image in another sequence. We then learn 7 models respectively for each cluster. Given a testing sequence, we choose the model trained from most similar cluster for prediction.

## Large-scale Labeling

We let the user label the 3D points in the 3D space. Using labels of 3D points, we are able to segment the 2D images at the same time.

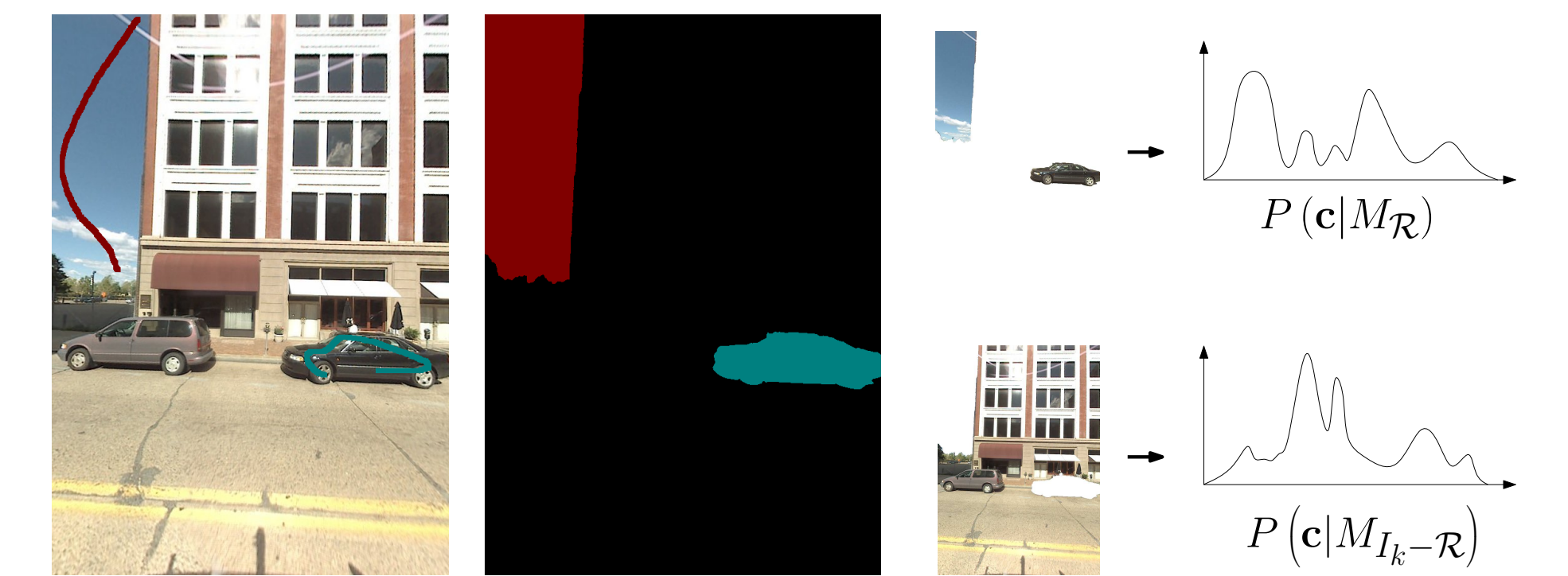


$$P_i^{3D}(l_i = l) \propto |\mathcal{T}^l \cap \mathcal{T}_i| + \frac{|\mathcal{T}^{\text{unknown}} \cap \mathcal{T}_i|}{n},$$

where  $\mathcal{T}^l = \{\mathbf{t} | \mathbf{t} \text{ is a track labeled as class } l \text{ by the user in 3D}\}$ ,  $\mathcal{T}^{\text{unknown}}$  is the set of feature tracks that have no label information from the user. The unary potential is defined to be

$$\psi_i(l) = -\frac{|\mathcal{T}_i - \mathcal{T}^{\text{unknown}}| + \epsilon}{H(P_i(\cdot)) + \epsilon} \log P_i(l),$$

where  $|\mathcal{T}_i - \mathcal{T}^{\text{unknown}}|$  is the number of labeled feature tracks with projections in superpixel  $p_i$ ,  $H(P_i(\cdot))$  is the entropy of the distribution  $P_i(\cdot)$ . For region with few 3D points, such as the sky, the user can do labeling in 2D.



When a superpixel in one image is covered by the strokes drawn by the user to be class  $l$ , the corresponding unary potential is set to  $\psi_i(l_i = l) = -\infty$  and  $\psi_i(l_i \neq l) = +\infty$ . For a superpixel without strokes, we determine the likelihood that the label information comes from the 2D color from

$$P_i^{\text{col}}(\mathbf{c}_i) = \frac{P(\mathbf{c}_i | M_{\mathcal{R}})}{P(\mathbf{c}_i | M_{\mathcal{R}}) + P(\mathbf{c}_i | M_{I_k - \mathcal{R}})}.$$

The probability is defined accordingly by

$$P_i(l) = P_i^{\text{col}}(\mathbf{c}_i) P_i^{2D}(l) + (1 - P_i^{\text{col}}(\mathbf{c}_i)) P_i^{3D}(l),$$

where  $P_i^{2D}(l)$  is the color likelihood computed from pixel colors in the stroke-covered regions of image belonging to class  $l$ .

## Results

