



Semantic Segmentation of Indoor-Scene RGB-D Images Based on Iterative Contraction and Merging

Jia-Hao Syu^{1(✉)}, Shih-Hsuan Cho², Sheng-Jyh Wang²,
and Li-Chun Wang¹

¹ Department of Communications Engineering, National Chiao Tung University,
Hsinchu, Taiwan
jiahaushiu4@gmail.com

² Department of Electronics Engineering, National Chiao Tung University,
Hsinchu, Taiwan

Abstract. In this paper, we propose an iterative contraction and merging framework (ICM) for semantic segmentation in indoor scenes. Given an input image and a raw depth image, we first derive the dense prediction map from a convolutional neural network (CNN) and a normal vector map from the depth image. By combining the RGB-D image with these two maps, we can guide the ICM process to produce a more accurate hierarchical segmentation tree in a bottom-up manner. After that, based on the hierarchical segmentation tree, we design a decision process which uses the dense prediction map as a reference to make the final decision of semantic segmentation. Experimental results show that the proposed method can generate much more accurate object boundaries if compared to the state-of-the-art methods.

Keywords: Convolutional neural network · Iterative contraction and merging
RGB-D image · Semantic segmentation

1 Introduction

Semantic segmentation, an important topic in computer vision, aims to assign each pixel a semantic label in an input image and to generate a dense semantic prediction for the given image. Up to now, many semantic segmentation algorithms have been proposed to improve the quality of dense semantic prediction. However, semantic segmentation is still a challenging work because of the complex and diverse contents in an indoor scene.

Today, RGB-D cameras are getting more popular and cheaper, such as Microsoft Kinect and Intel RealSense cameras. With RGB-D cameras, semantic segmentation algorithms [1–3] take into account both color and depth data to improve the quality of semantic labeling. On the other hand, deep learning techniques are getting popular due to the availability of large-scale datasets and powerful hardware. In [4], Long et al. proposed a deep learning model, called fully convolutional network (FCN), with both color and depth data to perform impressive semantic segmentation. However, since the

FCN model derives the dense semantic prediction by combining the up-sampling of multilayer information, the obtained semantic prediction results are usually not very accurate around the boundary area.

On the other hand, some hierarchical segmentation algorithms [5, 6] have adopted the bottom-up graph-based approach to generate a hierarchical segmentation tree for image segmentation. These algorithms can properly partition an image into image segments with very accurate region boundaries. However, due to the lack of semantic information during the bottom-up process, these hierarchical segmentation algorithms have difficulties in obtaining reasonable semantic segmentation results.

In this paper, we propose a semantic segmentation based on an iterative contraction and merging process to achieve semantic segmentation with much improved boundary extraction. In the proposed approach, we first cascade two bilateral filters to improve the quality of the raw depth data. Second, we integrate the RGB-D image with a dense semantic predictor, which extracts high-level information, and a normal estimation map, which extracts mid-level information, to guide the ICM process for the generation of a more accurate hierarchical segmentation tree. Finally, we make decisions over the hierarchical segmentation tree to obtain the final semantic segmentation result.

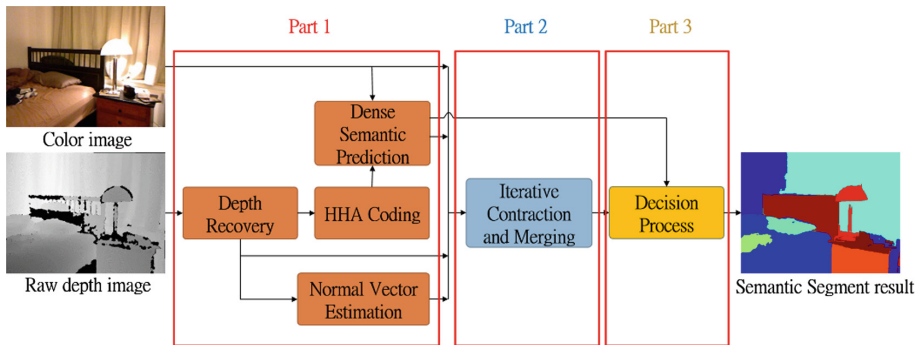


Fig. 1. Block diagrams of the proposed method

2 Proposed Method

In this paper, we propose an architecture that includes three parts to solve the semantic segmentation problem, as illustrated in Fig. 1. In the first part, we capture the mid-level information from the depth image, such as the refined depth image and the normal vector map, and the high-level information from the deep learning model, such as the dense semantic prediction map. To get the refined depth image, we cascade two bilateral filters to fill in the holes in the original depth image. Based on the refined depth image, we estimate the normal vector map based on cross-product computation. Moreover, the depth image is transformed into an HHA image, which has been defined in [3]. The FCN model combines both the RGB image and the HHA image to derive a semantic prediction map. In the second part, the iterative contraction and merging (ICM) process, which was originally proposed in [5], is used for unsupervised image

segmentation. In this paper, we maintain the original ICM features in [5] and add a few more features captured from the first part to derive a more robust hierarchical segmentation tree. In the final part, we design a decision process to decide the semantic segmentation result based on the hierarchical segmentation tree.

3 First Part: Feature Extraction

3.1 Depth Recovery

The raw depth data as shown in Fig. 2 (b) includes holes with no depth values. In order to fill in the estimated depth values over the holes, we assume adjacent pixels with similar RGB color values should have similar depth values. Besides, it is reasonable to trust nearby depth values than far-away depth values in the spatial domain. Hence, we design a spatial kernel which refers to the depth value. On the other hand, we design a range kernel which refers to the RGB values. Based on the spatial kernel and the range kernel, we define a bilateral filter $D^h(x)$ to fill in the depth values over the hole regions, as expressed below:

$$D^h(x) = W_1^{-1}(x) \int_{-k_1}^{k_1} \int_{-k_1}^{k_1} D(\xi) e^{-\frac{1}{2} \left(\frac{\|x-\xi\|}{\sigma_{s1}} \right)^2} e^{-\frac{1}{2} \left(\frac{\|I(\xi)-I(x)\|}{\sigma_{r1}} \right)^2} d\xi, \text{ with} \quad (1)$$

$$W_1(x) = \int_{-k_1}^{k_1} \int_{-k_1}^{k_1} e^{-\frac{1}{2} \left(\frac{\|x-\xi\|}{\sigma_{s1}} \right)^2} e^{-\frac{1}{2} \left(\frac{\|I(\xi)-I(x)\|}{\sigma_{r1}} \right)^2} d\xi, \quad (2)$$

where I and D denote the RGB and depth value, k_1 denotes the half-width of the filter, σ_{s1} denote the sigma of the spatial kernel, and σ_{r1} denotes the sigma of the range kernel. In our experiments, we choose $k_1 = 20$, $\sigma_{s1} = 10$ and $\sigma_{r1} = 0.1$. In some case, the hole area is too large in the original depth images so we need to iteratively fill the holes based on (1) until there is no hole in the depth image. In order to reduce the noise effect on the depth image, we apply the following bilateral filter to smooth the depth image:

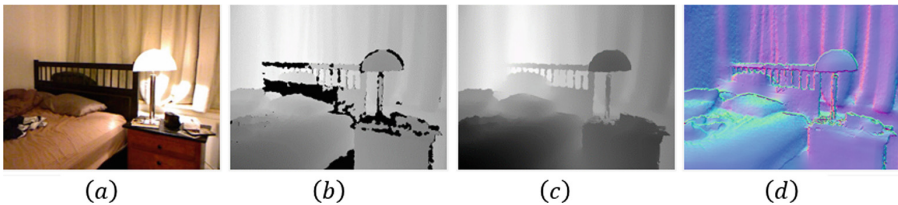


Fig. 2. (a) Original RGB image. (b) Original depth image. (c) Refined depth image. (d) Normal vector map.

$$D^s(x) = W_2^{-1}(x) \int_{-k_2}^{k_2} \int_{-k_2}^{k_2} D(\xi) e^{-\frac{1}{2} \left(\frac{\|x-\xi\|}{\sigma_{s2}} \right)^2} e^{-\frac{1}{2} \left(\frac{\|D(\xi)-D(x)\|}{\sigma_{r2}} \right)^2} d\xi, \text{ with} \quad (3)$$

$$W_2(x) = \int_{-k_2}^{k_2} \int_{-k_2}^{k_2} e^{-\frac{1}{2} \left(\frac{\|x-\xi\|}{\sigma_{s2}} \right)^2} e^{-\frac{1}{2} \left(\frac{\|D(\xi)-D(x)\|}{\sigma_{r2}} \right)^2} d\xi. \quad (4)$$

In our experiments, we choose $k_2 = 20$, $\sigma_{s2} = 12$ and $\sigma_{r2} = 0.05$. In Fig. 2(c), we show an example of the refined depth image.

3.2 Normal Vector Map

On the surface of the objects, each point is described as $(x, y, d(x, y))$ in the 3D coordinate where $d(x, y)$ denotes the depth value at (x, y) . We then estimate the gradient maps by computing the derivatives along the x-direction and y-direction of the depth values. On the surface S , the normal vector at a point i is derived by computing the cross product $\frac{\partial S_i}{\partial x} \times \frac{\partial S_i}{\partial y}$ which is expressed as

$$N_i = \frac{\partial S_i}{\partial x} \times \frac{\partial S_i}{\partial y} = \frac{\partial(x, y, d(x, y))}{\partial x} \times \frac{\partial(x, y, d(x, y))}{\partial y} = \left(-\frac{\partial d(x, y)}{\partial x}, -\frac{\partial d(x, y)}{\partial y}, 1 \right), \quad (5)$$

where $\frac{\partial d(x, y)}{\partial x} \approx \frac{d(x+1, y) - d(x-1, y)}{2}$ and $\frac{\partial d(x, y)}{\partial y} \approx \frac{d(x, y+1) - d(x, y-1)}{2}$. Based on the above computations, we can get the normal vector map as shown in Fig. 2(d).

3.3 Dense Semantic Prediction

In our approach, we generate a dense prediction map by using the fully convolutional network (FCN) [4]. The FCN model was originally designed for RGB images. In order to fit the 3-dimensional input format of the FCN model, we follow Gupta's approach in [3] and transform the refined depth image to the HHA image format. Besides, we fine-tune the FCN model for color images and learn another FCN model for depth images based on the NYUD-V2 dataset. Since the NYU-Depth V2 dataset [2] contains 40-class labels, the learned FCN model generates 40 layers of score maps. We combine each end of the feature map with the weight 0.5 and perform up-sampling to derive the score maps with each layer of the score maps representing one class. The dense prediction map is assigned by finding the maximal value at each pixel among 40 layers of score maps:

$$Class_i = \arg \min_K (S_{iK}), \quad (6)$$

where S_{iK} is the score of pixel i correspond to Class K .

4 Second Part: Iterative Contraction and Merging

The iterative contraction and merging process (ICM) in [5] can construct a hierarchical segmentation tree in two phases. In this paper, we aim to derive a more robust hierarchical segmentation tree for indoor-scene images. Hence, we maintain the original features captured in [5] and add additional features extracted in the first part of our system. In Phase 1 of the ICM process, the pixel-wise contraction and merging process quickly merges pixels with similar features into regions. In this phase, the definition of the affinity value $A(i, j)$ between pixels i and j is based on pixel-wise features, such as color, depth, normal vector, and dense semantic prediction value. After that, a remnant removal process is used to remove small remnant regions around the boundary. In Phase 2 of the ICM process, since similar image pixels have already been merged into regions, several kinds of regional information, including features in [5] and other additional features, such as depth, normal vector, and dense semantic prediction, are taken into account to define a more informative affinity value $A(R_m, R_n)$ to describe the similarity between the region pair R_m and R_n . Based on the region-wise affinity matrix, we iteratively apply the contraction and merging process to merge image regions into larger ones and progressively build a hierarchical segmentation tree. In the following subsections, we will describe the details of each module in the ICM process.

4.1 Phase-1: Pixel-Wise Contraction and Merging

Phase-1 of the ICM process aims to quickly merge pixels with similar features into regions. Here, we use a mixed feature space that consists of five subspaces: color space, spatial location space, depth space, normal vector space and dense prediction score space. The features at an image pixel i is mapped into a vector $(L_i, a_i, b_i, x_i, y_i, d_i, u_i, v_i, z_i, S_{i1} S_{i2} \dots S_{iK})$ in the feature space where (L_i, a_i, b_i) , (x_i, y_i) , (d_i) , (u_i, v_i, z_i) and $(S_{i1} S_{i2} \dots S_{iK})$ denote the color values, spatial coordinates, depth value, normal vector values and dense prediction values, respectively.

The contraction process aims to pull pixel pairs with similar features closer in the feature space than pixel pairs with less similar features. Here, the contraction process is formulated as the problem of finding the twisted coordinates $(\tilde{L}_i, \tilde{a}_i, \tilde{b}_i, \tilde{x}_i, \tilde{y}_i, \tilde{d}_i, \tilde{u}_i, \tilde{v}_i, \tilde{z}_i, \tilde{S}_{i1} \tilde{S}_{i2} \dots \tilde{S}_{iK})$ and is defined as

$$E(\tilde{\theta}) = \sum_{i=1}^N \sum_{j \in w_i, i \neq j} A(i, j) (\tilde{\theta}_i - \tilde{\theta}_j)^2 + \lambda_0 \sum_{i=1}^N (\tilde{\theta}_i - \theta_i)^2, \quad (7)$$

where N denotes the total number of pixels, w_i denotes the neighborhood around the pixel i , $\theta_i \in \{L_i, a_i, b_i, x_i, y_i, d_i, u_i, v_i, z_i, S_{i1} S_{i2} \dots S_{iK}\}$ denotes the original image feature values, $\tilde{\theta}_i \in \{\tilde{L}_i, \tilde{a}_i, \tilde{b}_i, \tilde{x}_i, \tilde{y}_i, \tilde{d}_i, \tilde{u}_i, \tilde{v}_i, \tilde{z}_i, \tilde{S}_{i1} \tilde{S}_{i2} \dots \tilde{S}_{iK}\}$ denotes the twisted image feature values, and $\tilde{\theta} = [\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_N]^T$. In our simulation, we empirically choose $\lambda_x = \lambda_y = 0.001$ and $\lambda_0 = 0.01$ otherwise. More details of contraction process can be found in [5]. In this paper, we define the affinity value $A(i, j)$ of pairwise pixels as

$$A(i, j) = \exp\left(-D(i, j) / \rho\right) \quad (8)$$

where $D(i, j)$ denotes the feature difference between i and j and is defined as

$$\begin{aligned} D(i, j) = & \kappa_1 \left\| [L_i \ a_i \ b_i]^T - [L_j \ a_j \ b_j]^T \right\|_2 + \kappa_2 \left\| [d_i]^T - [d_j]^T \right\|_2 + \kappa_3 \left\| [u_i \ v_i \ z_i]^T - [u_j \ v_j \ z_j]^T \right\|_2 \\ & + \kappa_4 \left\| [S_{i1} \ S_{i2} \dots \ S_{iK}]^T - [S_{j1} \ S_{j2} \dots \ S_{jK}]^T \right\|_2. \end{aligned} \quad (9)$$

Different from [5], here we add depth, normal vector, and dense prediction information into the affinity function. The score weight κ_1 , κ_2 , κ_3 and κ_4 controls the strength of the impact from the color, depth, normal vector and dense prediction, respectively. Similar to [5], the parameter ρ is adjusted to satisfy the condition that 70% of the $A(i, j)$ values is larger than 0.01.

After applying the contraction process, we use the same grid-based merging process in [5] to group nearby pixels into regions. In order to efficiently perform the merging process, we only consider the $(L_i, a_i, b_i, x_i, y_i)$ features during merging. Here, the feature space $S^\theta = \max(\theta) - \min(\theta)$ denotes the dynamic range in each feature value. To achieve grid-based merging, we divided the feature space from $S^L \times S^a \times S^b \times S^x \times S^y$ into $[S^L/15][S^a/15][S^b/15][S^x/25][S^y/25]$ regions.

After the pixel-wise contraction and merging process, pixels are merged into regions. However, there may exist a few pixels around the boundary of objects looking like noisy data. To deal with this problem, the remnant removal process in [5] is used to merge regions whose size is smaller than the predefined threshold into one of their adjacent regions with the most similar color appearance.

4.2 Phase-2: Region-Wise Contraction and Merging

After Phase-1, image pixels have been merged into a set of regions. Similar to Phase-1, the averaged feature values in each region R_m can be mapped into a feature vector $(L_m, a_m, b_m, x_m, y_m, d_m, u_m, v_m, z_m, S_{m1}S_{m2} \dots S_{mK})$ in the feature space. Likewise, we can derive the twisted coordinates and define the energy function as

$$E(\tilde{\theta}) = \sum_{m=1}^{N_R} \sum_{R_n \in \phi} A(R_m, R_n) \left(\tilde{\theta}_{R_m} - \tilde{\theta}_{R_n} \right)^2 + \lambda_{\theta} \sum_{m=1}^{N_R} \left(\tilde{\theta}_{R_m} - \theta_{R_m} \right)^2 \quad (10)$$

where N_R denotes the total number of regions, ϕ denotes the neighboring regions of R_m , $\theta_{R_m} \in \{L_m, a_m, b_m, x_m, y_m, d_m, u_m, v_m, z_m, S_{m1}S_{m2} \dots S_{mK}\}$ denotes the original feature values, $\tilde{\theta}_{R_m} \in \{\tilde{L}_m, \tilde{a}_m, \tilde{b}_m, \tilde{x}_m, \tilde{y}_m, \tilde{d}_m, \tilde{u}_m, \tilde{v}_m, \tilde{z}_m, \tilde{S}_{m1}\tilde{S}_{m2} \dots \tilde{S}_{mK}\}$ denotes the twisted feature values, and $\tilde{\theta} = [\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_{N_R}]^T$. The affinity function of pairwise region is defined as

$$A(R_m, R_n) = \exp\left(-D(R_m, R_n)/\rho\right), \quad (11)$$

where the parameter ρ is set so that 10% of the affinity values are larger than 0.01 if $N_R > 200$ and ρ is set so that 1% of the affinity values are larger than 0.01 if $N_R \leq 200$. The difference function between two regions R_m and R_n is defined as

$$D(R_m, R_n) = D_N(R_m, R_n) \left[M_1 \left(\frac{\alpha D_C(R_m, R_n) + \beta D_T(R_m, R_n) + \gamma \tilde{D}_C^B(R_m, R_n)}{\sqrt{(80 + SI(R_m, R_n))}} \right) \right. \\ \left. + (M_2 D_D(R_m, R_n) + M_3 \tilde{D}_D^B(R_m, R_n)) + M_4 D_N(R_m, R_n) + M_5 D_S(R_m, R_n) \right]. \quad (12)$$

where the terms D_N , D_C , D_T , \tilde{D}_C^B and SI denote the difference of region-size, color, texture, color-of-border and spatial-intertwining. Some of these terms have been introduced in [5]. In this paper, we add four additional difference terms, including depth D_D , depth-of-border \tilde{D}_D^B , normal-vector D_N and dense prediction score D_S . In our simulation, we empirically choose $\alpha = 1$, $\beta = 3$, $\gamma = 6$, $M_1 = 3$, $M_2 = 1$, $M_3 = 6$, $M_4 = 1$ and $M_5 = 2$. After the ICM, we can construct a hierarchical segmentation tree. Please refer to [5] for the details of hierarchical segmentation tree construction.

5 Decision Process

Based on the hierarchical segmentation tree, we propose a decision process to find a semantic segmentation by referring to the dense prediction map. In the dense prediction map, we merge pixels with the same class into regions. For each region S_K in the dense semantic prediction map, we check the corresponding node in the hierarchical segmentation tree that has the largest overlapping with S_K . We call this region the candidate region of S_K . In other words, given a class region S_K from the dense semantic prediction map and a node region T_n , the candidate region C_K is defined as

$$C_K = \arg \min_n \left(\frac{|S_K \cap T_n|}{|T_n|}, \frac{|T_n|}{|S_K \cap T_n|} \right), \quad (13)$$

where n denotes the node in the hierarchical segmentation tree and $|\cdot|$ denotes number of pixels in the region. After computing candidate regions, we calculate the covering in the image with these candidate regions. Here, three cases may occur at each pixel:

- (1) one candidate region: the pixel is only covered by one candidate region,
- (2) more than one candidate region: the pixel is covered by more than two candidate regions, and
- (3) no candidate region: the pixel is not covered by any candidate region.

For the first case, we can immediately assign the semantic label based on the corresponding class label. In the second case, we tend to trust the candidate region with a smaller size and assign the semantic label accordingly. In the third case, we have to

assign these no-candidate regions with some semantic labels. Here, we reverse the search from the no-candidate region into the dense prediction map. We can search multiple nodes in the no-candidate region and find the larger overlap with the class of the dense prediction map.

6 Experimental Results

Our propose model is evaluated on NYU-Depth V2 dataset [2] which includes 1449 RGBD images captured by Microsoft Kinect V1. The dataset contains dense per-pixel labeling and are classified into 40 class for semantic segmentation task by Gupta et al. [3]. The quantitative evaluation is measured by four common metrics: pixel accuracy, mean accuracy, mean IU, and frequency weighted IU. In order to know the most important features in our proposed method, we add depth, normal vector, and dense prediction score. In Table 1, we list the quantitative evaluation over the 100 random sampling images in NYUD-V2 dataset by using different combinations of the RGB, depth, normal vector, and dense prediction score. It can be observed that the high-level dense prediction score provides more important information than other features. It turns out the combination of all features does provide the most preferred results. In Table 2, we also compare our proposed method with other semantic segmentation methods. We can observe that the quantitative evaluation of our proposed method is only close to the original FCN model. This is because most boundary regions of the objects are not taken into account in the quantitative evaluation and the improvement of our method mainly occur in those boundary areas.

Table 1. Quantitative evaluation over different combinations of distance functions. (1) RGB, (2) Depth, (3) Normal Vector and (4) FCN

ICM + DP	Pixel acc.	Mean acc.	Mean IU	f.w. IU
1	64.9	51.3	37.9	51.7
1, 2	65.9	51.8	38.6	52.5
1, 3	66.0	52.3	39.0	53.0
1, 4	66.9	53.8	40.1	54.7
1, 2, 3, 4	67.5	54.2	41.0	54.9

Table 2. Comparison of semantic segmentation methods over 100 random sampling images on NYUD-V2 dataset (average value of 5 experiments)

Architecture	Pixel acc.	Mean acc.	Mean IU	f.w. IU
Gupta et al. [3]	60.3	28.6	31.3	47.0
FCN [4]	67.9	56.6	42.4	55.0
Ours	67.5	54.2	41.0	54.9

In Fig. 3, we compare the semantic segmentation results of our method with the original FCN [4] method over four image samples. It can be easily seen that our proposed method provides more accurate semantic segmentation results around the object boundaries.

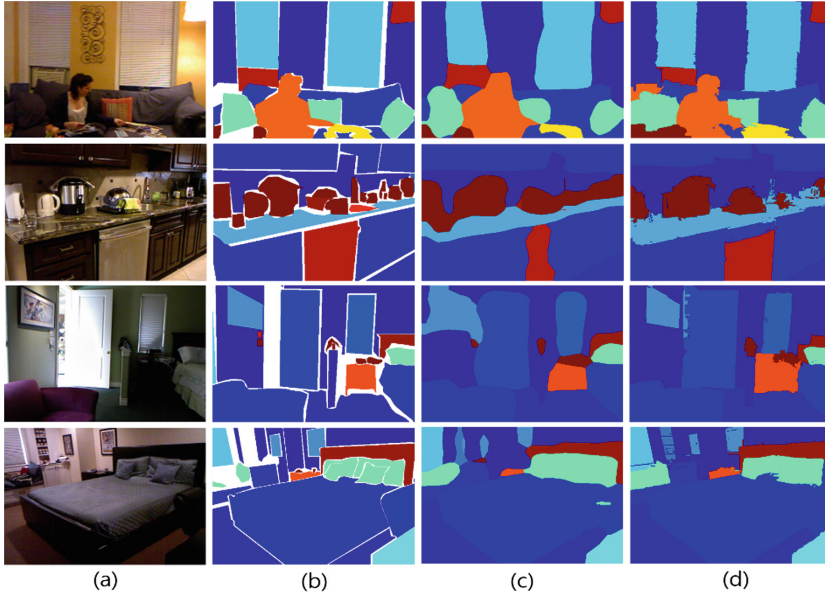


Fig. 3. (a) Original image. (b) Ground truth. (c) FCN [4] method. (d) Our proposed method.

7 Conclusion

In this paper, we propose an iterative contraction and merging framework for semantic segmentation of indoor-scene images. Based on the ICM framework, we improve the quality of the hierarchical segmentation tree by considering more mid-level and high-level features. We also design a decision process to decide the final semantic segmentation result based on the hierarchical segmentation tree. Experimental results show that the proposed method can generate more accurate object boundaries on semantic segmentation results.

References

1. Ren, X., Bo, L., Fox, D.: RGB-(D) scene labeling: features and algorithms. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 2759–2766 (2012)
2. Silberman, N., Hoiem, D., Kohli, P., Fergus, R.: Indoor segmentation and support inference from RGBD images. In: Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., Schmid, C. (eds.) ECCV 2012 Part V. LNCS, vol. 7576, pp. 746–760. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-33715-4_54

3. Gupta, S., Girshick, R., Arbeláez, P., Malik, J.: Learning rich features from RGB-D images for object detection and segmentation. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) ECCV 2014 Part VII. LNCS, vol. 8695, pp. 345–360. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-10584-0_23
4. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 3431–3440 (2015)
5. Syu, J.H., Wang, S.J., Wang, L.C.: Hierarchical image segmentation based on iterative contraction and merging. *IEEE Trans. Image Process.* **26**(5), 2246–2260 (2017)
6. Arbelaez, P., Maire, M., Fowlkes, C., Malik, J.: Contour detection and hierarchical image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **33**(5), 898–916 (2011)