

A Coarse-to-Fine Indoor Layout Estimation (CFILE) Method

Yuzhuo Ren, Chen Chen, Shangwen Li, and C.-C. Jay Kuo

Abstract. The task of estimating the spatial layout of cluttered indoor scenes from a single RGB image is addressed in this work. Existing solutions to this problems largely rely on hand-craft features and vanishing lines, and they often fail in highly cluttered indoor rooms. The proposed coarse-to-fine indoor layout estimation (CFILE) method consists of two stages: 1) coarse layout estimation; and 2) fine layout localization. In the first stage, we adopt a fully convolutional neural network (FCN) to obtain a coarse-scale room layout estimate that is close to the ground truth globally. The proposed FCN considers combines the layout contour property and the surface property so as to provide a robust estimate in the presence of cluttered objects. In the second stage, we formulate an optimization framework that enforces several constraints such as layout contour straightness, surface smoothness and geometric constraints for layout detail refinement. Our proposed system offers the state-of-the-art performance on two commonly used benchmark datasets.

1 Introduction

The task of spatial layout estimation of indoor scenes is to locate the boundaries of the floor, walls and the ceiling. It is equivalent to the problem of semantic surface labeling. The segmented boundaries and surfaces are valuable for a wide range of computer vision applications such as indoor navigation [1], object detection [2] and augmented reality [1,3,4,5]. Estimating the room layout from a single RGB image is a challenging task. This is especially true in highly cluttered rooms since the ground and wall boundaries are often occluded by various objects. Besides, indoor scene images can be shot at different viewpoints with large intra-class variation. As a result, high-level reasoning is often required to avoid confusion and uncertainty. For example, the global room model and its associated geometric reasoning can be exploited for this purpose. Some researchers approach this layout problem by adding the depth information [6,7].

The indoor room layout estimation problem has been actively studied in recent years. Hedau *et al.* [8] formulated it as a structured learning problem. It first generates hundreds of layout proposals based on inference from vanishing lines. Then, it uses the line membership features and the geometric context features to rank the obtained proposals and chooses the one with the highest score as the desired final result.

In this work, we propose a coarse-to-fine indoor layout estimation (CFILE) method. Its pipeline is shown in Fig. 1. The system uses an RGB image as its

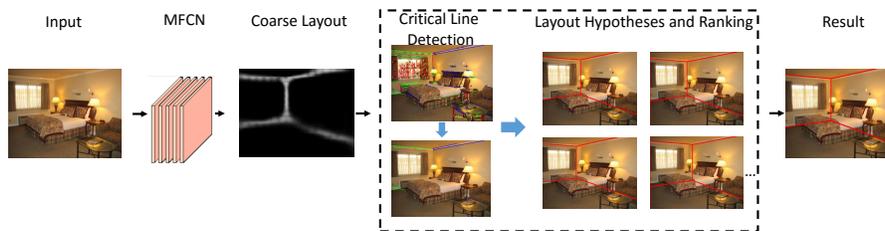


Fig. 1. The pipeline of the proposed coarse-to-fine indoor layout estimation (CFILe) method. For an input indoor image, a coarse layout estimate that contains large surfaces and their boundaries is obtained by a multi-task fully convolutional neural network (MFCN) in the first stage. Then, occluded lines and missing lines are filled in and possible layout choices are ranked according to a pre-defined score function in the second stage. The one with the highest score is chosen to the final output.

input and provides a box layout as its output. The CFILe method consists of two stages: 1) coarse layout estimation; and 2) fine layout localization. In the first stage, we adopt a multi-task fully convolutional neural network (MFCN) [9] to obtain a coarse-scale room layout estimate. This is motivated by the strength of the FCN in semantic segmentation [10] and contour detection [11]. The FCN has a strong discriminant power in handling a large variety of indoor scenes using the surface property and the layout contour property. It can provide a robust estimate in the presence of cluttered objects, which is close to the ground truth globally. In the second stage, being motivated by structured learning, we formulate an optimization framework that enforces several constraints such as layout contour straightness, surface smoothness and geometric constraints for layout detail refinement.

It is worthwhile to emphasize that the spatial layout estimation problem is different from semantic object segmentation problem in two aspects. First, the spatial layout problem targets at the labeling of semantic surface of an indoor room rather than objects in the room. Second, we have to label occluded surfaces while semantic segmentation does not deal with the occlusion problem at all. It is also different from the contour detection problem since occluded layout contours have to be detected.

The major contributions of this work are three folds. First, we use the FCN to learn the labeling of main surfaces and key contours jointly, which are critical to robust spatial layout of an indoor scene. The FCN training is elaborated. It is shown that the course-scale layout estimate obtained by the FCN is robust and close to ground truth. Second, we formulate an optimization framework that enforces three constraints (i.e. surface smoothness, contour straightness and proper geometrical structure) to refine the coarse-scale layout estimate. Third, we conduct extensive performance evaluation by comparing the proposed CFILe

method and several benchmarking methods on the dataset of Hedau *et al.* [8], the LSUN validation dataset [12]. It is shown by experimental results that the proposed CFILE method offers the state-of-the-art performance. It outperforms the second best method by 1.16% and 1.32% in Hedau’s dataset and the LSUN dataset, respectively.

The rest of this paper is organized as follows. Related previous work is reviewed in Sec. 2. The proposed CFILE method is described in detail in Sec. 3. Experimental results are shown in Sec. 4. Concluding remarks are drawn in Sec. 5.

2 Related Work

Structured Learning. The structured learning methodology [13] has been widely used in the context of indoor room layout estimation. It targets at learning the structure of an environment in the presence of imperfect low-level features. It consists of two stages [13]. First, a set of structure hypotheses are generated. Second, a score function is defined to evaluate the structure in hypotheses set. The first stage is guided by low level features such as vanishing lines under the Manhattan assumption. The number of layout hypotheses in the first stage is usually large while most of them are of low accuracy due to the presence of clutters. If the quality of hypotheses is low in the first stage, there is no easy way to fix it in the second stage. In the second stage of layout ranking, the score function contains various features such as the line membership [8,14], the geometric context [8,14], the object location [15], etc. The score function cannot handle objects well since they overlap with more than one surfaces (e.g., between the floor and walls). The occluding objects in turn make the surface appearance quite similar along their boundaries.

Classical Methods for Indoor Layout Estimation. Research on indoor room layout estimation has been active in recent years. Hedau *et al.* [8] formulated it as a structured learning problem. There are many follow-up efforts after this milestone work. They focus on either developing new criteria to reject invalid layout hypotheses or introducing new features to improve the score function in layout ranking.

Different hypothesis evaluation methods were considered in [7,8,15,16,17,18,19]. Hedau *et al.* [8] reduced noisy lines by removing clutters first. Specifically, they used the line membership together with semantic labeling to evaluate hypotheses. Gupta *et al.* [15] proposed an orientation map that labels three orthogonal surface directions based on line segments and, then, used the orientation map to re-evaluate layout proposals. Besides, they detected objects and fit them into 3D boxes. Since an object cannot penetrate the wall, they used the box location as a constraint to reject invalid layout proposals. The work in [2,20] attempted to model objects and spatial layout simultaneously. Hedau *et al.* [21] improved their earlier work in [2,8] by localizing the box more precisely using several cues such as edge- and corner-based features. Ramalingam *et al.* [19] proposed an algorithm to detect Manhattan Junctions and selected the best layout by

optimizing a conditional random field whose corners are well aligned with pre-detected Manhattan Junctions. Pero *et al.* [18] integrated the camera model, an enclosing room box, frames (windows, doors, pictures), and objects (beds, tables, couches, cabinets) to generate layout hypotheses. Lampert *et al.* [22] improved objects detection by maximizing a score function through the branch and bound algorithm.

3D- and Video-based Indoor Layout Estimation. Zhao and Zhu [17] exploited the location information and 3D spatial rules to obtain as many 3D boxes as possible. For example, if a bed is detected, the algorithm will search its neighborhood to look for a side table. Then, they rejected impossible layout hypothesis. Choi *et al.* [23] trained several 3D scene graph models to learn the relation among the scene type, the object type, the object location and layout jointly. Guo *et al.* [7] recovered 3D model from a single RGBD image by transferring the exemplar layout in the training set to the test image. Fidler *et al.* [24] and Xiang *et al.* [25] represented objects by a deformable 3D cuboid model for improved object detection and then used in layout estimation. Fouhey *et al.* [26] exploited human action and location in time-lapse video to infer functional room geometry.

CNN- and FCN-based Indoor Layout Estimation. The convolution neural network (CNN) has a great impact on various computer vision research topics, such as object detection, scene classification, semantic segmentation, etc. Mallya and Lazebnik [14] used the FCN to learn the informative edge from an RGB image to provide a rough layout. The FCN shares features in convolution layers and optimize edges detection and geometric context labeling [8,27,28] jointly. The learned contours are used as a new feature in sampling vanishing lines for layout hypotheses generation. Dasgupta *et al.* [29] used the FCN to learn semantic surface labels. Instead of learning edges, their solution adopted the heat map of semantic surfaces obtained by the FCN as the belief map and optimized it furthermore by vanishing lines. Generally speaking, a good layout should satisfy several constraints such as boundary straightness, surface smoothness and proper geometrical structure. However, the CNN is weak in imposing spatial constraints and performing spatial inference. As a result, an inference model was appended in both [14] and [29] to refine the layout result obtained by CNN.

3 Coarse-to-Fine Indoor Layout Estimation (CFILE)

3.1 System Overview

Most research on indoor layout estimation [7,8,15,16,17,18,19] is based on the “Manhattan World” assumption. That is, a room contains three orthogonal directions indicated by three groups of vanishing lines. Hedau *et al.* [8] presented a layout model based on 4 rays and a vanishing point. The model can written as

$$\text{Layout} = (l_1, l_2, l_3, l_4, v), \quad (1)$$

where l_i is the i^{th} line and v is the vanishing point. If (l_1, l_2, l_3, l_4, v) can be easily detected without any ambiguity, the layout problem is straightforward.

One example is given in Fig. 2 (a), where five surfaces are visible in the image without occlusion.

However, more challenging cases exist. Vertices p_i and e_i in Fig. 2 (a) may lie outside the image. One example is shown in Fig. 2 (b). Furthermore, vertices p_2 and p_3 are floor corners and they are likely to be occluded by objects. Furthermore, line l_2 may be entirely or partially occluded as shown in Fig. 2 (c). Lines l_3 and l_4 are wall boundaries, and they can be partially occluded but not fully occluded. Line l_1 is the ceiling boundary which is likely to be visible.

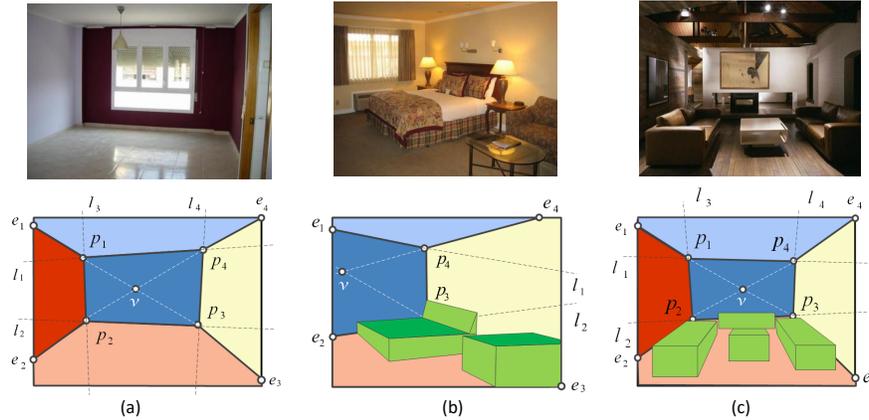


Fig. 2. Illustration of a layout model $\text{Layout} = (l_1, l_2, l_3, l_4, v)$ that is parameterized by four lines and a vanishing point: (a) an easy setting where all five surfaces are present; (b) a setting where some surfaces are outside the image; (c) a setting where key boundaries are occluded.

The proposed CFILe system consists of two stages as illustrated in Fig. 1. In the first stage, we propose a multi-task fully convolutional neural network (MFCN) to offer a coarse yet robust layout estimation. Since the CNN is weak in imposing spatial smoothness and conducting geometric reasoning, it cannot provide a fine-scale layout result. In the second stage, we first use the coarse layout from MFCN as the guidance to detect a set of critical lines. Then, we generate a small set of high quality layout hypotheses based on these critical lines. Finally, we define a score function to select the best layout as the desired output. Detailed tasks in these two stages are elaborated below.

3.2 Coarse Layout Estimation via MFCN

We adopt a multi-task fully convolutional neural network (MFCN) [10,14,9] to learn the coarse layout of indoor scenes. The MFCN [9] shares features in the convolutional layers with those in the fully connected layers and builds different

branches for multi-task learning. The total loss of the MFCN is the sum of losses of different tasks. The proposed two-task network structure is shown in Fig. 3. We use the VGG-16 architecture for fully convolutional layers and train the MFCN for two tasks jointly, i.e. one for layout learning while the other for semantic surface learning (including the floor, left-, right-, center-walls and the ceiling). Our work is different from that in [14], where layout is trained together with geometric context labels [27,28] which contains object labels. Here, we train the layout and semantic surface labels jointly. By removing objects from the concern, the boundaries of semantic surfaces and layout contours can be matched even in occluded regions, leading to a clearer layout. As compared to the work in [29], which adopts the fully convolutional neural network to learn semantic surfaces with a single task network, our network has two branches, and their learned results can help each other.

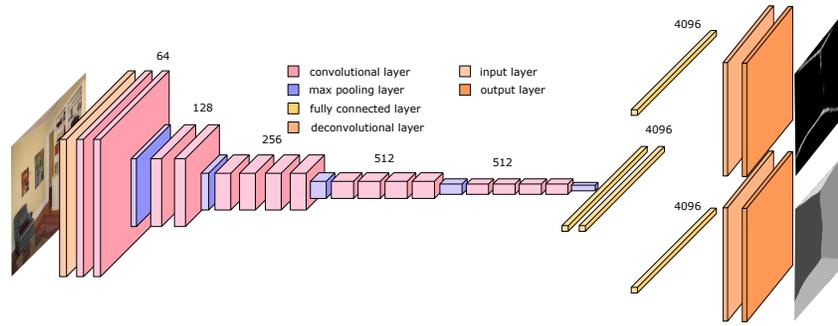


Fig. 3. Illustration of the FCN-VGG16 with two output branches. We use one branch for the coarse layout learning and the other branch for semantic surface learning. The input image size is re-sized to 404×404 to match the receptive field size of the filter at the fully connection layer.

The receptive field of the filter at the fully connected layer of the FCN-VGG16 is 404×404 , which is independent of the input image size [10,30]. Xu *et al.* [30] attempted to vary the FCN training image size so as to capture different level of details in image content. If the input image size is larger than the receptive field size, the filter of the fully connected layer looks at a part of the image. If the input image size is smaller than the receptive field size, it is padded with zeros and spatial resolution is lost in this case. The layout describes the whole image’s global structure. We resize the input image to 404×404 so that the filter examines the whole image.

3.3 Layout Refinement

There are two steps in structured learning: 1) to generate a hypotheses set; and 2) to define a score function and search a structure in the hypotheses set that maximizes the score function. We attempt to improve in both areas.

Given an input image \mathbf{I} of size $w \times h \times 3$, the output of the coarse layout from the proposed MFCN in Fig. 3 is a probability function in form of

$$\mathbf{P}^{(k)} = Pr(\mathbf{L}_{ij} = k|\mathbf{I}), \quad \forall k \in \{0, 1\}, i \in [1, \dots, h], j \in [1, \dots, w], \quad (2)$$

where \mathbf{L} is an image of size $w \times h$ that maps each pixel in the original image, \mathbf{I}_{ij} , to a label in the output image $\mathbf{L}_{ij} \in \{0, 1\}$, where 0 denotes a background pixel and 1 denotes a layout pixel. One way to estimate the final layout from the MFCN output is to select the label with the highest score; namely,

$$\hat{\mathbf{L}}_{ij} = \underset{k}{\operatorname{argmax}} \mathbf{P}_{ij}^{(k)} \quad \forall i \in [1, \dots, h], j \in [1, \dots, w]. \quad (3)$$

It is worthwhile to point out that $\hat{\mathbf{L}}_{ij}$ generated from the MFCN output is noisy for two reasons. First, the contour from the MFCN is thick and not straight since the convolution operation and the pooling operation lose the spatial resolution gradually along stages. Second, the occluded floor boundary (e.g., the l_2 line in Fig. 2) is more difficult to detect since it is less visible than other contours (e.g., the l_1 , l_3 and l_4 lines in Fig. 2). We need to address these two challenges in defining a score function.

The optimal solution for Eq. (3) is difficult to get directly. Instead, we first generate layout hypotheses that are close to the global optimal layout, denoted by \mathbf{L}^* , in the layout refinement algorithm. Then, we define a novel score function to rank layout hypotheses and select the one with the highest score as the final result.

Generation of High-Quality Layout Hypotheses Our objective is to find a set of layout hypotheses that contains fewer yet more robust proposals in the presence of occluders. Then, the best layout with the smallest error can be selected.

Vanishing Line Sampling. We first threshold the layout contour obtained by the MFCN, convert it into a binary mask, and dilate it by 4 pixels to get a binary mask image denoted by C . Then, we apply the vanishing lines detection algorithm [15] to the original image and select those inside the binary mask as critical lines $l_{i(\text{original})}$, shown in solid lines in Fig. 4 (c) (d) (e) for ceiling, wall and floor separately. Candidate vanishing point v is generated by grid search around the initial v from [15].

Handling Undetected Lines. There is case when no vanishing lines are detected inside C because of low contrast, such as wall boundaries, l_3 (or l_4). If ceiling corners are available, l_3 (or l_4) are filled in by connecting ceiling corners and vertical vanishing point. If ceiling corners do not present in the image, the missing l_3 (or l_4) is estimated by logistic regression use the layout points in \mathbf{L} .

Handling Occluded Lines. As discussed earlier, the floor line, l_2 , can be entirely or partially occluded. One illustrative example is shown in Fig. 4 where l_2 is partially occluded. If l_2 is partially occluded, the occluded part of l_2 can be recovered by line extension. For entirely occluded l_2 , if we simply search lines inside C or uniformly sample lines [14], the layout proposal is not going to be accurate as the occluded boundary line cannot be recovered. Instead, we automatically fill in occluded lines based on geometric rule. If p_2 (or p_3) is detectable by connecting detected l_3 (or l_4) to e_2v (or e_3v), l_2 is computed as the line passing through the available p_2 or p_3 and the vanishing point l_2 associated with. If neither p_2 nor p_3 is detectable, l_2 is estimated by logistic regression use the layout points in \mathbf{L} .

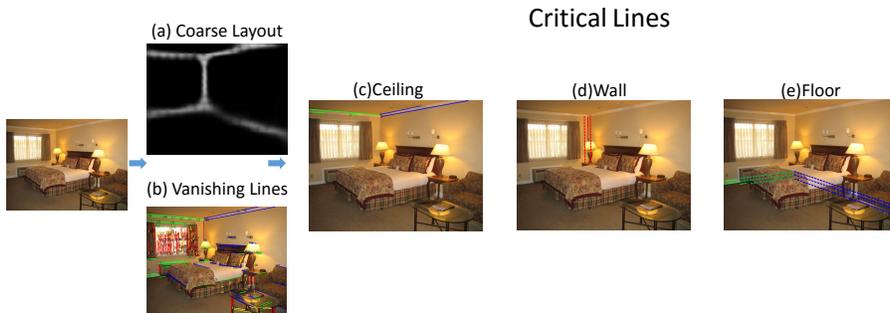


Fig. 4. Illustration of critical lines detection for better layout hypothesis generation. For a given input image, the coarse layout offers a mask that guides vanishing lines selection and critical lines inference. The solid lines indicate detected vanishing lines C . The dashed wall lines indicate those wall lines that are not detected but inferred inside mask C from ceiling corners. The dashed floor lines indicate those floor lines that are not detected but inferred inside mask C .

In summary, the final l_{critical} used in generating layout hypotheses is the union of three parts as given below:

$$l_{\text{critical}} = l_{i(\text{original})} \cup l_{i(\text{occluded})} \cup l_{i(\text{undetected})}, \quad (4)$$

where $l_{i(\text{original})}$ denotes detected vanishing lines inside C , $l_{i(\text{occluded})}$ denotes the recovered occluded boundary, and $l_{i(\text{undetected})}$ denotes undetected vanishing lines because of low contrast but recovered from geometric reasoning. These three types of lines are shown in Fig. 4. With $l_{i(\text{original})}$ and vanishing point v , we generate all possible layouts \mathbf{L} using the model described in Sec. 3.1.

Layout Ranking We use the coarse layout probability map \mathbf{P} as a weight mask to evaluate the layout. The score function is defined as

$$S(\mathbf{L}|\mathbf{P}) = \frac{1}{N} \sum_{i,j} \mathbf{P}_{i,j}, \quad \forall \mathbf{L}_{i,j} = 1, \quad (5)$$

where \mathbf{P} is the output from the MFCN, \mathbf{L} is a layout from the hypotheses set, N is a normalization factor that is equal to the total number of layout pixels in \mathbf{L} . Then, the optimal layout is selected by

$$\mathbf{L}^* = \underset{\mathbf{L}}{\operatorname{argmax}} S(\mathbf{L}|\mathbf{P}). \quad (6)$$

The score function is in favor of the layout that is aligned well with the coarse layout. Fig. 5 shows one example where the layout hypotheses are ranked using the score function in Eq. (6). The layout with the highest score is chosen to be the final result.

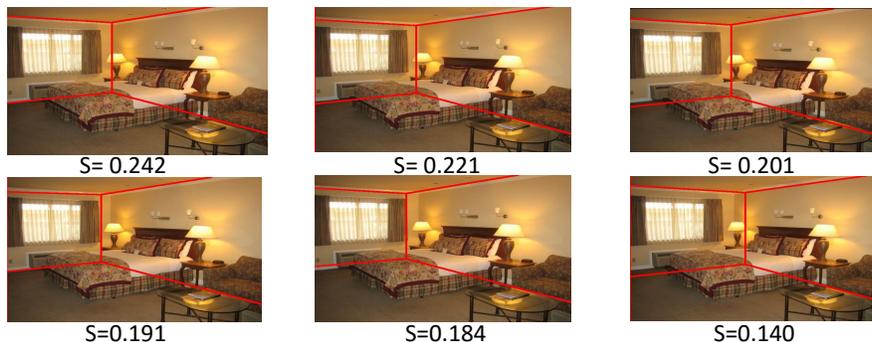


Fig. 5. Example of Layout ranking using the proposed score function.

4 Experiments

4.1 Experimental Setup

We evaluate the proposed CFILe method on two popular datasets; namely, Hedau’s dataset [8] and the LSUN dataset [14]. Hedau dataset contains 209 training images, 53 validation images and 105 test images. Mallya *et al.* [14] expanded Hedau dataset by adding 75 new images into training set while validation and test set unchanged, which referred to Hedau+ dataset. We conduct data augmentation for Hedau+ dataset as done in [14] by cropping, rotation, scaling and luminance adjustment in the training of the MFCN. The LSUN dataset [14] contains 4000 training images, 394 validation images and 1000 test images. Since

no ground truth is released for the 1000 test images, we evaluate the proposed method on the 394 validation set only. We resize all images to 404×404 by bicubic interpolation in the MFCN training, and train two coarse layout models for the two datasets separately.

Hedau+ dataset provides both the layout and the geometric context labels but it does not provide semantic surface labels. Thus, we use the layout polygon provided in the dataset to generate semantic surface labels. The LSUN dataset provides semantic surface labels but not the layout. We detect edges on semantic surface labels and dilate them to a width of 7 pixels in the MFCN training. By following [14], we use the NYUDv2 RGBD dataset in [31] for semantic segmentation to initialize the MFCN. Also, we set the base learning rate to 10^{-4} with momentum 0.99.

We adopt two performance metrics: the pixel-wise error and the corner error. To compute the pixel-wise error, the obtained layout segmentation is mapped to the ground truth layout segmentation. Then, the pixel-wise error is the percentage of pixels that are wrongly matched. To compute the corner error, we sum up all Euclidean distances between obtained corners and their associated ground truth corners.

4.2 Experimental Results and Discussion

The coarse layout scheme described in Sec. 3.2 is first evaluated using the methodology in [32]. We compare our results, denoted by MFCN_1 and MFCN_2 , against the informative edge method [14], denoted by FCN, in Table 1. Our proposed two coarse layout schemes have higher ODS (fixed contour threshold) and OIS (per-image best threshold) scores. This indicates that they provide more accurate regions for vanishing line samples in layout hypotheses generation.

Table 1. Performance comparison of coarse layout results for Hedau’s test dataset, where the performance metrics are the fixed contour threshold (ODS) and the per-image best threshold (OIS) [32]. We use FCN to indicate the informative edge method in [14]. Both MFCN_1 and MFCN_2 are proposed in our work. They correspond to the two settings where the layout and semantic surfaces are jointly trained on the original image size (MFCN_1) and the downsampled image size 404×404 . (MFCN_2)

Metrics	FCN[14]		$\text{MFCN}_1(\text{our})$		$\text{MFCN}_2(\text{our})$	
	ODS	OIS	ODS	OIS	ODS	OIS
Hedau’s dataset	0.255	0.263	0.265	0.284	0.265	0.291

We use several exemplary images to demonstrate that the proposed coarse layout results are robust and close to the ground truth. That is, we compare visual results of the FCN in [14] and the proposed MFCN_2 in Fig. 6. As compared to the layout results of the FCN in [14], the proposed MFCN_2 method provides

robust and clearer layout results in occluded regions, which are not much affected by object boundaries.

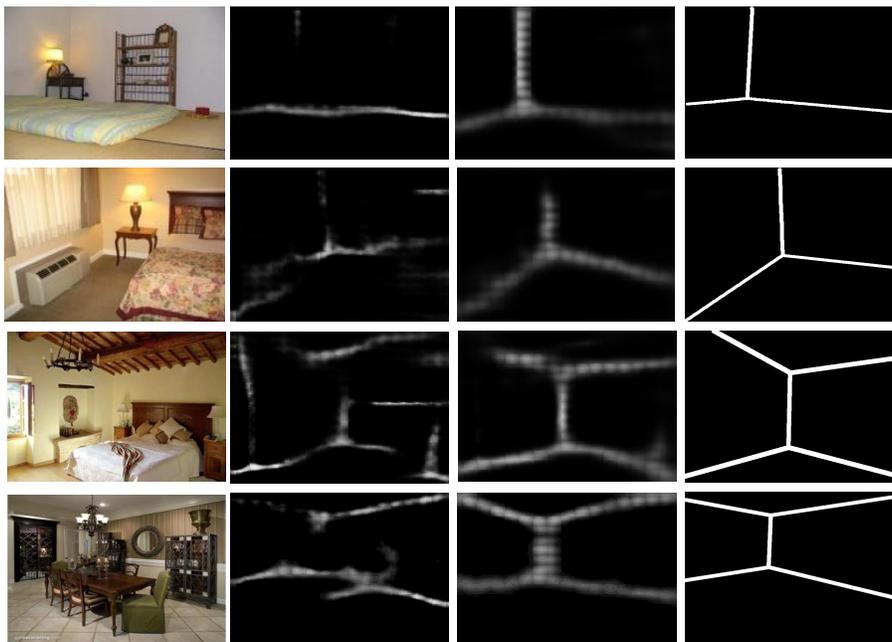


Fig. 6. Comparison of coarse layout results (from left to right): the input image, the coarse layout result of the FCN in [14], the coarse layout results of the proposed $MFCN_2$ and the ground truth. The results of the $MFCN_2$ are more robust. Besides, it provides clearer contours in occluded regions. The first two examples are from Hedau dataset and the last two examples are from LSUN dataset.

Next, we evaluate the performance of the proposed full layout algorithm, CFILE, including the coarse layout estimation and the layout optimization and ranking. The performance of several methods for Hedau’s dataset and the LSUN dataset is compared in Table 2 and Table 3, respectively. The proposed CFILE method achieves state-of-the-art performance. It outperforms the second best algorithm by 1.16% in Hedau’s dataset and 1.32% in the LSUN dataset.

The best six results of the proposed CFILE method for Hedau’s test images are visualized in Fig. 7. We see from these five examples that the coarse layout estimation algorithm is robust in highly cluttered rooms (see the second row and the fourth). The layout refinement algorithm can recover occluded boundaries accurately in Fig. 7 (a), (b), (d) and (e). It can also select the best layout among several possible layouts. The worst three results of the proposed CFILE method

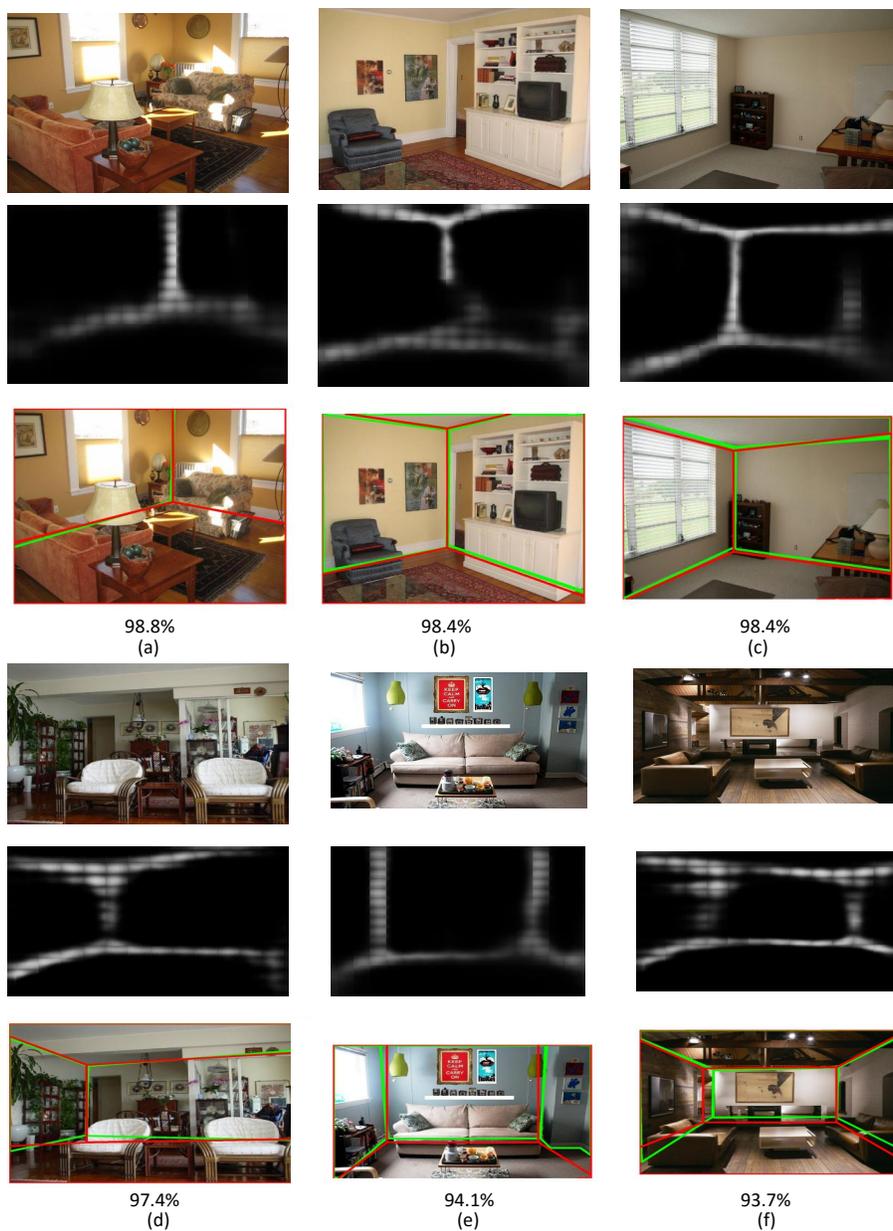


Fig. 7. Visualization of six best results of the CFILE method in Hedau's test dataset (from top to bottom): original images, the coarse layout estimates from MFCN, our results with pixel-wise accuracy (where the ground truth is shown in green and our result is shown in red).

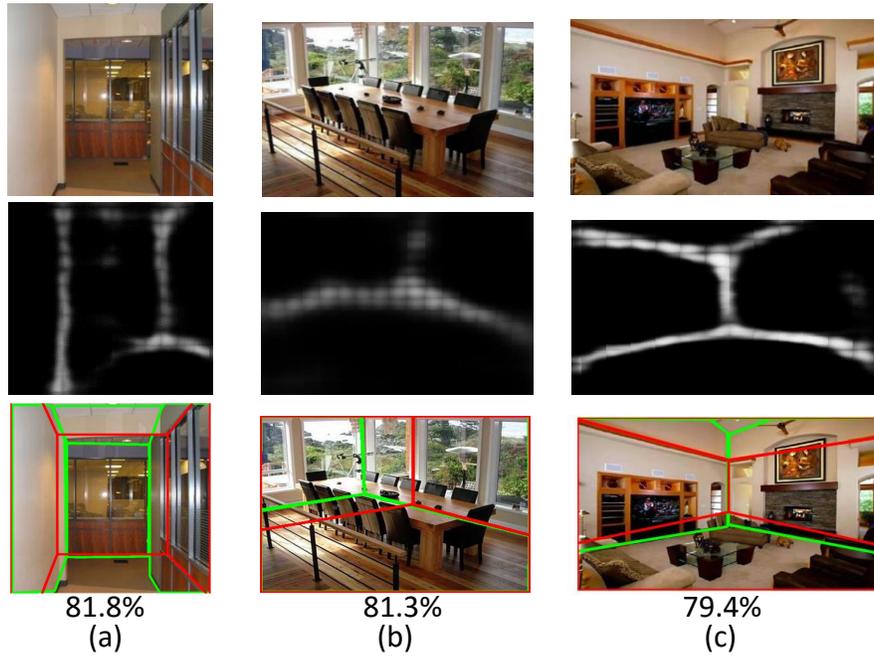


Fig. 8. Visualization of three worst results of the CFILE method in Hedau’s test dataset (from top to bottom): original images, the coarse layout estimates from MFCN, our results with pixel-wise accuracy (where the ground truth is shown in green and our result is shown in red).

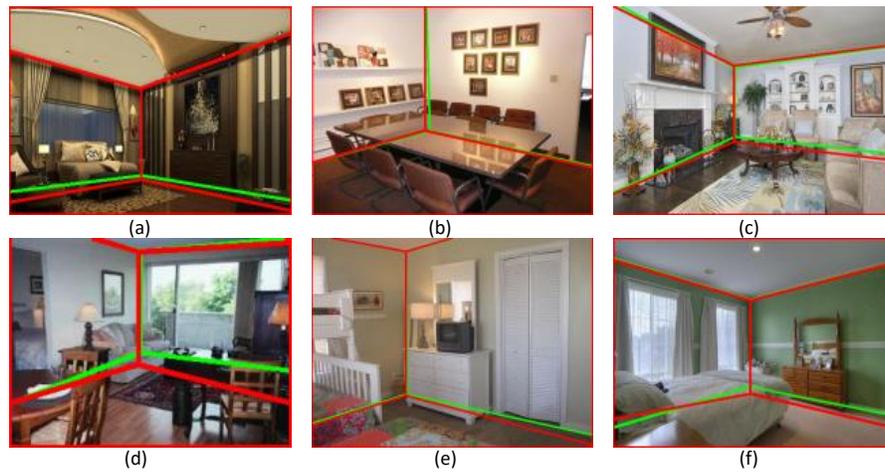


Fig. 9. Visualization of layout results of the CFILE method in the LSUN validation set. Ground truth is shown in green and our result is shown in red.

Table 2. Performance benchmarking for Hedau’s dataset.

Method	Pixel Error (%)
Hedau <i>et al.</i> (2009)[8]	21.20
Del Pero <i>et al.</i> (2012)[18]	16.30
Gupta <i>et al.</i> (2010)[15]	16.20
Zhao <i>et al.</i> (2013)[17]	14.50
Ramalingam <i>et al.</i> (2013)[19]	13.34
Mallya <i>et al.</i> (2015)[14]	12.83
Schwing <i>et al.</i> (2012)[33]	12.80
Del Pero <i>et al.</i> (2013)[34]	12.70
Dasgupta <i>et al.</i> (2016)[29]	9.73
Proposed CFILE	8.67

Table 3. Performance benchmarking for the LSUN dataset.

Method	Corner Error (%)	Pixel Error (%)
Hedau <i>et al.</i> (2009)[8]	15.48	24.23
Mallya <i>et al.</i> (2015)[14]	11.02	16.71
Dasgupta <i>et al.</i> (2016) [29]	8.20	10.63
Proposed CFILE	7.95	9.31

for Hedau’s test images are visualized in Fig. 8. Fig. 8 (a) show one example where the fine layout result is misled by the wrong coarse layout estimate. Fig. 8 (b) is a difficult case. The left wall and right wall have the same appearance and there are several confusing wall boundaries. Fig. 8 (c) gives the worst example of the CFILE method with accuracy 79.4%. However, it is still higher than the worst example reported in [14] with accuracy 61.05%. The ceiling boundary is confusing in Fig. 8 (f). The proposed CFILE method selects the ceiling line overlapping with the coarse layout. More visual results from the LSUN dataset are shown in Fig. 9.

5 Conclusion and Future Work

A coarse-to-fine indoor layout estimation (CFILE) method was proposed to estimate the room layout from an RGB image. We adopted a multi-task fully convolutional neural network (MFCN) to offer a robust coarse layout estimate for a variety of indoor scenes with joint layout and semantic surface training. However, CNN is weak in enforcing spatial constraints. To address this problem, we formulated an optimization framework that enforces several constraints such as layout contour straightness, surface smoothness and geometric constraints for layout detail refinement. It was demonstrated by experimental results that the proposed CFILE system offers the best performance on two commonly used benchmark datasets. It is an interesting topic to investigate the multi-scale effect

of CNN-based vision solutions and their applications to semantic segmentation and geometrical layout of indoor scenes.

References

1. Karsch, K., Hedau, V., Forsyth, D., Hoiem, D.: Rendering synthetic objects into legacy photographs. In: *ACM Transactions on Graphics (TOG)*. Volume 30., ACM (2011) 157
2. Hedau, V., Hoiem, D., Forsyth, D.: Thinking inside the box: Using appearance models and context based on room geometry. In: *Computer Vision–ECCV 2010*. Springer (2010) 224–237
3. Liu, C., Schwing, A.G., Kundu, K., Urtasun, R., Fidler, S.: Rent3d: Floor-plan priors for monocular layout estimation. In: *Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on, IEEE (2015) 3413–3421*
4. Xiao, J., Furukawa, Y.: Reconstructing the worlds museums. *International Journal of Computer Vision* **110** (2014) 243–258
5. Martin-Brualla, R., He, Y., Russell, B.C., Seitz, S.M.: The 3d jigsaw puzzle: Mapping large indoor spaces. In: *Computer Vision–ECCV 2014*. Springer (2014) 1–16
6. Zhang, J., Kan, C., Schwing, A., Urtasun, R.: Estimating the 3d layout of indoor scenes and its clutter from depth sensors. In: *Proceedings of the IEEE International Conference on Computer Vision*. (2013) 1273–1280
7. Guo, R., Zou, C., Hoiem, D.: Predicting complete 3d models of indoor scenes. *arXiv preprint arXiv:1504.02437* (2015)
8. Hedau, V., Hoiem, D., Forsyth, D.: Recovering the spatial layout of cluttered rooms. In: *Computer vision, 2009 IEEE 12th international conference on, IEEE (2009) 1849–1856*
9. Dai, J., He, K., Sun, J.: Instance-aware semantic segmentation via multi-task network cascades. *arXiv preprint arXiv:1512.04412* (2015)
10. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2015) 3431–3440
11. Xie, S., Tu, Z.: Holistically-nested edge detection. In: *Proceedings of the IEEE International Conference on Computer Vision*. (2015) 1395–1403
12. <http://lsun.cs.princeton.edu/2015.html>.
13. Nowozin, S., Lampert, C.H.: Structured learning and prediction in computer vision. *Foundations and Trends® in Computer Graphics and Vision* **6** (2011) 185–365
14. Mallya, A., Lazebnik, S.: Learning informative edge maps for indoor scene layout prediction. In: *Proceedings of the IEEE International Conference on Computer Vision*. (2015) 936–944
15. Gupta, A., Hebert, M., Kanade, T., Blei, D.M.: Estimating spatial layout of rooms using volumetric reasoning about objects and surfaces. In: *Advances in neural information processing systems*. (2010) 1288–1296
16. Schwing, A., Fidler, S., Pollefeys, M., Urtasun, R.: Box in the box: Joint 3d layout and object reasoning from single images. In: *Proceedings of the IEEE International Conference on Computer Vision*. (2013) 353–360
17. Zhao, Y., Zhu, S.C.: Scene parsing by integrating function, geometry and appearance models. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2013) 3119–3126

18. Pero, L.D., Bowdish, J., Fried, D., Kermgard, B., Hartley, E., Barnard, K.: Bayesian geometric modeling of indoor scenes. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, IEEE (2012) 2719–2726
19. Ramalingam, S., Pillai, J., Jain, A., Taguchi, Y.: Manhattan junction catalogue for spatial reasoning of indoor scenes. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2013) 3065–3072
20. Wang, H., Gould, S., Roller, D.: Discriminative learning with latent variables for cluttered indoor scene understanding. *Communications of the ACM* **56** (2013) 92–99
21. Hedau, V., Hoiem, D., Forsyth, D.: Recovering free space of indoor scenes from a single image. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, IEEE (2012) 2807–2814
22. Lampert, C.H., Blaschko, M.B., Hofmann, T.: Efficient subwindow search: A branch and bound framework for object localization. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **31** (2009) 2129–2142
23. Choi, W., Chao, Y.W., Pantofaru, C., Savarese, S.: understanding indoor scenes using 3d geometric phrases. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2013) 33–40
24. Fidler, S., Dickinson, S., Urtasun, R.: 3d object detection and viewpoint estimation with a deformable 3d cuboid model. In: *Advances in Neural Information Processing Systems*. (2012) 611–619
25. Xiang, Y., Savarese, S.: Estimating the aspect layout of object categories. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, IEEE (2012) 3410–3417
26. Fouhey, D.F., Delaitre, V., Gupta, A., Efros, A.A., Laptev, I., Sivic, J.: People watching: Human actions as a cue for single view geometry. *International Journal of Computer Vision* **110** (2014) 259–274
27. Hoiem, D., Efros, A.A., Hebert, M.: Geometric context from a single image. In: *Computer Vision, 2005. ICCV 2005. Tenth IEEE International Conference on*. Volume 1., IEEE (2005) 654–661
28. Hoiem, D., Efros, A.A., Hebert, M.: Recovering surface layout from an image. *International Journal of Computer Vision* **75** (2007) 151–172
29. Saumitro Dasgupta, Kuan Fang, K.C.S.S.: Delay: Robust spatial layout estimation for cluttered indoor scenes. In: *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on*, IEEE (2016)
30. Xu, H., Venugopalan, S., Ramanishka, V., Rohrbach, M., Saenko, K.: A multi-scale multiple instance video description network. *arXiv preprint arXiv:1505.05914* (2015)
31. Gupta, S., Arbelaez, P., Malik, J.: Perceptual organization and recognition of indoor scenes from rgb-d images. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2013) 564–571
32. Arbelaez, P., Maire, M., Fowlkes, C., Malik, J.: Contour detection and hierarchical image segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **33** (2011) 898–916
33. Schwing, A.G., Urtasun, R.: Efficient exact inference for 3d indoor scene understanding. In: *Computer Vision—ECCV 2012*. Springer (2012) 299–313
34. Pero, L., Bowdish, J., Kermgard, B., Hartley, E., Barnard, K.: Understanding bayesian rooms using composite 3d object models. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2013) 153–160