# SHALLOW2DEEP: INDOOR SCENE MODELING BY SINGLE IMAGE UNDERSTANDING

**Yinyu Nie**
Bournemouth University
Poole BH12 5BB, U.K.
ynie@bournemouth.ac.uk

**Shihui Guo**[*]
Xiamen University
Xiamen 361005, China.
guoshihui@xmu.edu.cn

**Jian Chang**[*]
Bournemouth University
Poole BH12 5BB, U.K.
jchang@bournemouth.ac.uk

**Xiaoguang Han**
The Chinese University of Hong Kong (Shenzhen)
Shenzhen 518172, China.
hanxiaoguang@cuhk.edu.cn

**Jiahui Huang**
Tsinghua University
Beijing 100084, China
huang-jh18@mails.tsinghua.edu.cn

**Shi-Min Hu**
Tsinghua University
Beijing 100084, China
shimin@tsinghua.edu.cn

**Jian Jun Zhang**
Bournemouth University
Poole BH12 5BB, U.K.
jzhang@bournemouth.ac.uk

## ABSTRACT

Dense indoor scene modeling from 2D images has been bottlenecked due to the absence of depth information and cluttered occlusions. We present an automatic indoor scene modeling approach using deep features from neural networks. Given a single RGB image, our method simultaneously recovers semantic contents, 3D geometry and object relationship by reasoning indoor environment context. Particularly, we design a shallow-to-deep architecture on the basis of convolutional networks for semantic scene understanding and modeling. It involves multi-level convolutional networks to parse indoor semantics/geometry into non-relational and relational knowledge. Non-relational knowledge extracted from shallow-end networks (e.g. room layout, object geometry) is fed forward into deeper levels to parse relational semantics (e.g. support relationship). A Relation Network is proposed to infer the support relationship between objects. All the structured semantics and geometry above are assembled to guide a global optimization for 3D scene modeling. Qualitative and quantitative analysis demonstrates the feasibility of our method in understanding and modeling semantics-enriched indoor scenes by evaluating the performance of reconstruction accuracy, computation performance and scene complexity.

*K*eywords Scene understanding · Image-based modeling · Semantic modeling · Relational reasoning

## 1 Introduction

Understanding indoor environment is of significant impact and has already been applied in the domains of interior design, real estate, etc. In recent years, 3D scanning and reconstruction of indoor scenes have been intensively explored using various sensors [1]. Understanding 3D indoor contents from an RGB image shows its unique significance for our daily-life applications, e.g. 3D digital content generation for social media and content synthesis for virtual reality and augmented reality.

3D scene modeling from a single image is challenging as it requires computers to perform equivalently as human vision to perceive and understand indoor context with only color intensities. It generally requires for blending various

---
[*]Corresponding authors

vision tasks [1] and most of them are still under active development, e.g. object segmentation [2], layout estimation [3] and geometric reasoning [4]. Although machine intelligence has reached comparable human-level performance in some tasks (e.g. scene recognition [5]), those techniques are only able to represent a fragment knowledge of full scene context.

With the lack of depth clues, prior studies reconstructed indoor scenes from a single image by exploiting shallow image features (e.g. line segments and HOG descriptors [6, 4]) or introducing depth estimation [7, 8] to search object models. Other works adopt Render-and-Match strategy to obtain CAD scenes with their renderings similar as input images [9]. However, it is still an unresolved problem when indoor geometry is over-cluttered and complicated. The reasons are threefold. First, complicated indoor scenes involve heavily occluded objects, which could cause missing contents in detection [9]. Second, cluttered environments significantly increase the difficulty of camera and layout estimations, which critically affects the reconstruction quality [10]. Third, compared to the large diversity of objects in real scenes, the reconstructed virtual environment is still far from satisfactory (missing small pieces, wrong labeling). Existing methods have explored the use of various contextual knowledge, including object support relationship [7, 8] and human activity [7], to improve modeling quality. However, their relational (or contextual) features are hand-crafted and would fail to cover a wide range of objects in cluttered scenes.
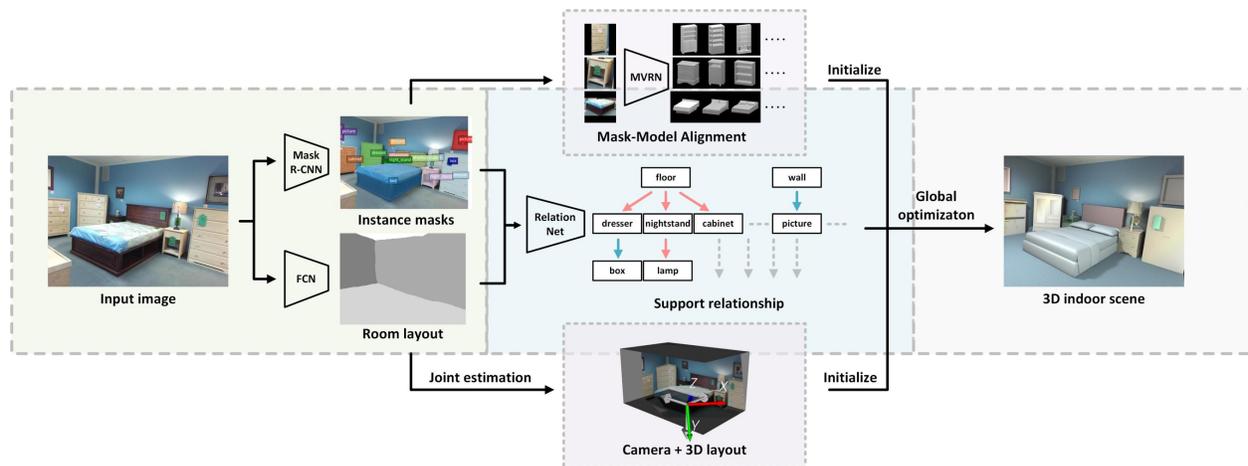


Figure 1: Pipeline of indoor scene modeling from a single image. The whole process is divided into three phases: 1. non-relational semantics parsing (e.g. room layout and object masks); 2. support relationship inference; 3. global scene optimization.

Different from previous works, our work aims at dense scene modeling. We extract and assemble object semantics (i.e. object masks with labels) and geometric contents (i.e. room layout and object models) into structured knowledge after being processed with shallower stacks of neural networks (see Figure 1). All the extracted semantics and geometry above is passed to deeper stacks of networks to infer support relationships between objects, which guides the final 3D scene modeling with global optimization. We take advantage of these object support relations and improve the modeling performance, in terms of object diversity and accuracy. We also propose a novel method to jointly estimate the camera parameters and room layout, which helps to improve scene modeling accuracy compared with using the existing methods. In summary, our contributions can be listed as follows:

- a Relation Network for support inference. This network predicts support relationship between indoor objects. It improves the reconstruction quality in the stage of global scene optimization, particularly increasing the accuracy of object placement in occlusions.

- a global optimization strategy for indoor scene synthesis. It incorporates the outputs from former networks and iteratively recovers 3D scenes to make them contextually consistent with the scene context. It performs effectively in inferring the shape of heavily occluded objects.

- a unified scene modeling system backboned by convolutional neural networks (CNN). With the capability of CNNs in parsing scene contents, latent image features are perceived and accumulated from consecutive networks. It outputs compact indoor context with a shallow-to-deep streamline to automatically generate semantics-enriched 3D scenes.

## 2    Related Work

**Indoor Content Parsing**    Capturing indoor contents is a prerequisite for modeling semantic scenes. Prior studies have explored scene parsing from various aspects with a single image input, which can be divided into two branches by their targets: 1) semantic detection and 2) geometric reasoning.

In semantic detection, deep learning holds satisfying performance in extracting latent features. It provides high accuracy in acquiring various types of scene semantics, e.g. scene types (bedroom, office or living room, etc.) [11], object labels (bed, lamp or picture etc.) [12], instance masks [13] and room layout (locations of wall, floor and ceiling) [10, 14]. On these raw semantics, abstract descriptions like scene grammar have been used to summarize them into a hierarchical structure. It divides indoor contents into groups by functionality or spatial relationship (e.g. affiliation relations) for scene understanding [15]. Unlike semantic detection, geometric reasoning refers to capturing spatial clues from images, e.g. depth map [16], surface normal map [17] and object geometry (including model retrieval [18] and reconstruction [19]).

However, these scene parsing methods are well tailored for a specific task. We design our holistic scene parsing step on their basis. Semantic and geometric clues are extracted and unified to form structured knowledge for scene understanding and 3D modeling.

**Support Inference**    This step is inspired by the research [20] where a freehand sketch drawing is turned into semantically valid, well arranged 3D scenes. They performed co-retrieval and co-placement of 3D models which are related to each other (e.g. a chair and a desk) by jointly processing the sketched objects in pair. The success of their work showed the significance of relations between objects, while our work contributes to this importance element by automatically inferring the object support relationship.

Support relationship provides a sort of geometric constraint between indoor objects to build scenes more robustly. This originates from daily experience that an object requires some support to counteract the gravity. Support inference from RGB images is an ambiguous problem without knowing the 3D geometry, where occlusions usually make the supporting parts invisible in the field of view. However, the arrangement of indoor furniture generally follows a set of interior design principles and living habits (e.g. tables are mostly supported by the floor; pictures are commonly on walls). These latent patterns behind scenes make the support relationship a kind of priors that can be learned by viewing various indoor rooms. Earlier studies addressed this problem by designing specific optimization problems [21, 22] considering both depth clues and image features. Apart from inferring support relations, many researchers represented the support condition by object stability and indoor safety [23, 24]. Moreover, the support relation is also a type of spatial relationship in scene grammar to enhance contextual bindings between objects. Other approaches implemented support inference to understand scenes with a scene graph [15, 7]. However, these methods either require for depth information, or rely on hand-crafted priors or models. In our work, we formulate the support inference problem into a Visual-Question Answering (VQA) form [25], where a Relation Network is designed to end-to-end link the object pairs and infer their support relations.

**Single-view Scene Modeling**    Indoor scene modeling from RGB images can be divided into two branches: 1. layout prediction and 2. indoor content modeling. Based on the Manhattan assumption [26], layout prediction represents indoor layout with cuboid proposals using line segments [27] or CNN feature maps [14, 10].

To reconstruct indoor contents, previous methods adopted cuboid shapes [28, 29] to recover the orientation and placement of target objects without the need for querying CAD model datasets. However, the geometric details are weak because objects are only represented by a bounding box. Rather than using cuboid shapes, other methods produced promising results in placement estimation of a single object by aligning CAD models with the object image [30, 31]. Other methods leveraged shallow features (e.g. line segments, edges and HOG features) [6, 32, 4] to segment images and retrieve object models, or used a scene dataset as priors to retrieve object locations based on co-occurrence statistics [33]. They either asked for human interaction or hand-crafted priors in parsing scene contents, while semantics and object geometry are learned with our method, allowing the capability of handling extendable object categories.

Recent studies also considered CNNs to detect objects and estimate room layout [9, 8] with informative scene knowledge [7]. Huang et al. [7] estimated depth maps and surface normal maps from RGB images with scene grammar to optimize the object placement. However, depth prediction is sensitive if the input distribution is slightly different from the training data [8]. Instead of tailoring scene grammar to improve the reconstruction results, we incorporates the relational reasoning in our process to infer the object relationship with a Relation Network. A parallel development [9] followed a Render-and-Match strategy to optimize object locations and orientations, which did not involve any depth clues and other relational constraints. CAD scenes are iterated until their renderings are sufficiently close to the input images. In contrast, our method does not refer to extra depth prediction and scene rendering iterations, which significantly boosts

the computing performance. The scene modeling is built on a unified CNN architecture. Intermediate semantic and geometric clues are parsed and accumulated by sub-modules, and reorganized with support relations to guide the scene modeling.

## 3    Overview

Our framework is built on the hypothesis that, features produced in each phase could be accumulated to feed into the consequent networks for deeper scene understanding. This process is divided into three phases, as illustrated in Figure 1. The first phase obtains non-relational semantics (i.e.room layout, object masks and labels) and retrieves a small set of 3D object candidates from a large model library (Section 4). This part takes advantage of a number of recent development in computer vision communities. We tailored a selection of methods to precondition the non-relational features for solving the scene modeling problem in later two phases.

In the second phase, we introduce a Relation Network to infer support relationships between objects (Section 5). This relational semantics offers physical constraints to organize those non-relational information into a reasonable contextual structure for 3D modeling.

The third phase assembles the geometric contents to model the 3D scene contextually consistent with these relational and non-relational semantics (Section 6). The 3D room layout and camera orientation are jointly estimated to ensure their consistence. It provides two coordinate systems (the room coordinate system and the camera coordinate system) for the global optimization in scene modeling and refinement.

## 4    Non-Relational Semantics Parsing

**2D Layout Estimation**    Layout estimation provides room boundary geometry (i.e. the location of the floor, ceiling and walls). Using CNNs to produce layout features, current works [14, 10] generally ask for camera parameters to estimate vanishing points for layout proposal decision. We adopt the Fully Convolutional Network (FCN) from [14] to extract the layout edge map and label map. These feature maps present a coarse prediction of 2D layout (see Appendix D). An accurate 3D layout is jointly estimated along with camera parameters for further scene modeling (see Section 6.1).

**Scene Segmentation**    We segment images at the instance-level to obtain object category labels and corresponding 2D masks. Object masks present meaningful clues to initialize object sizes, 3D locations and orientations. Particularly, we introduce the Mask R-CNN [13] to capture object masks with instance segmentation. We customize the Mask R-CNN backboned by ResNet-101 [34], with the weights pre-trained on the MSCOCO dataset [35]. It is fine-tuned on the NYU v2 dataset [21] which contains 1,449 densely labeled indoor images covering 37 common and 3 'other' categories. (The training strategy is detailed in Appendix A.1). Since object masks act significantly in the latter stages, we append Mask R-CNN with the Dense Conditional Random Field (DCRF) [36] to merge overlaps and improve mask edges. Besides, wall, floor and ceiling masks are removed as they are precisely decided in the 2D layout estimation. Segmentation samples are shown in Figure 2.

**Model Retrieval**    This task is to retrieve CAD models with the most similar appearance to the segmented object images. A Multi-View Residual Network (MVRN) pretrained on ShapeNet [37] is introduced for shape retrieval. Similar with [9, 7], we align and render each model from 32 viewpoints (two elevation angles at 15 and 30 degrees, and 16 uniform azimuth angles) for appearance-based matching. A Multi-View Convolutional Network [38] backboned with ResNet-50 is adopted as feature extractors to view CAD models from different viewpoints. This type of architecture is designed to mimic human eyes by observing an object from multiple viewpoints to recognize the object shape. Deep features from a single view is represented by a 2048-dimensional vector (i.e. the last layer size of ResNet-50). This compact descriptor enables us to match models efficiently in the vector space. The similarity between an image and a model can be measured by the cosine distance: $\max_{i\in[1,32]}\cos(\boldsymbol{f}, \boldsymbol{f}_i^{\mathrm{m}})$, $\boldsymbol{f}, \boldsymbol{f}_i^{\mathrm{m}} \in \mathcal{R}^{2048}$, where $\boldsymbol{f}$ and $\boldsymbol{f}_i^{\mathrm{m}}$ respectively denote the shape descriptor of the object image and a rendering of the model. The model set construction and training strategy are detailed in Appendix A.2. Furthermore, we fine-tune the orientation of matched models with ResNet-34 (see Discussions). Figure 3 shows some matched samples on our model set. Top-5 candidates are selected for global scene optimization in Section 6.

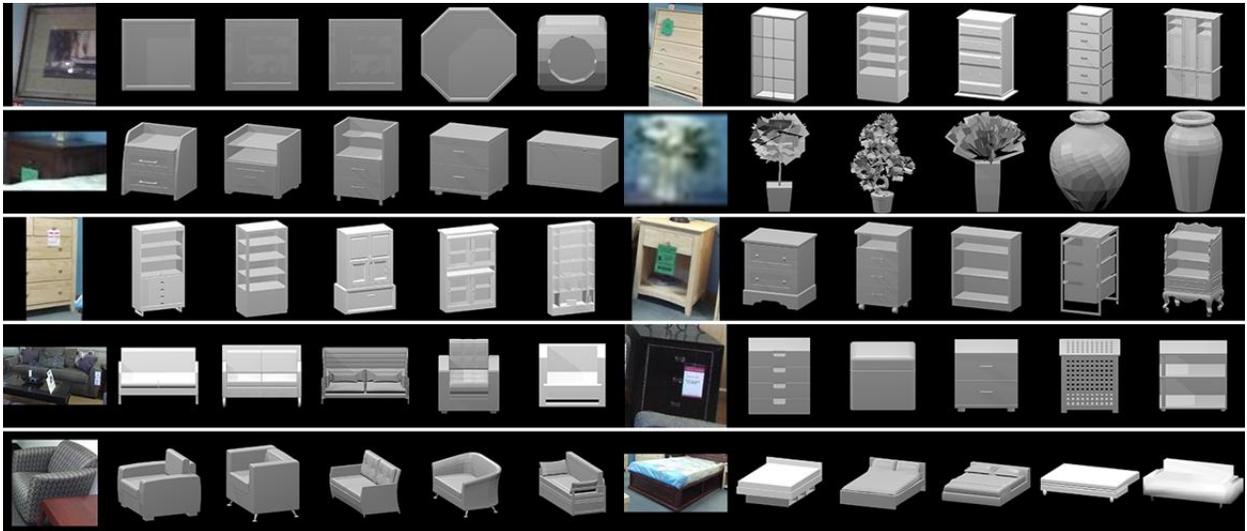Figure 2: Instance segmentation samples



Figure 3: CAD model candidates. For each object image, we search our model dataset and output five similar candidates for scene modeling.

## 5 Relational Reasoning

Section 4 dedicates to parsing indoor scenes into non-relational contents. We here aim to extract relational clues from these upstream outputs to conclude support relationships between objects. This relationship serves as physical constraints to guide scene modeling.

As assumed in existing works [21, 39], two support types are considered in our paper (i.e. support from behind, e.g. on a wall, or below, e.g. on a table). Every object except layout instances (i.e. wall, ceiling and floor) must be supported by another one. For objects which are supported by hidden instances, we treat them as being supported by layout instances.

Unlike non-relational semantics, relational context asks for not only the object property features, but also the contextual link between object pairs. Thus, a key is to combine the object feature pairs with specific task descriptions for support reasoning. It can be intuitively formulated as a Visual Question Answering (VQA) manner [25, 40]: given the
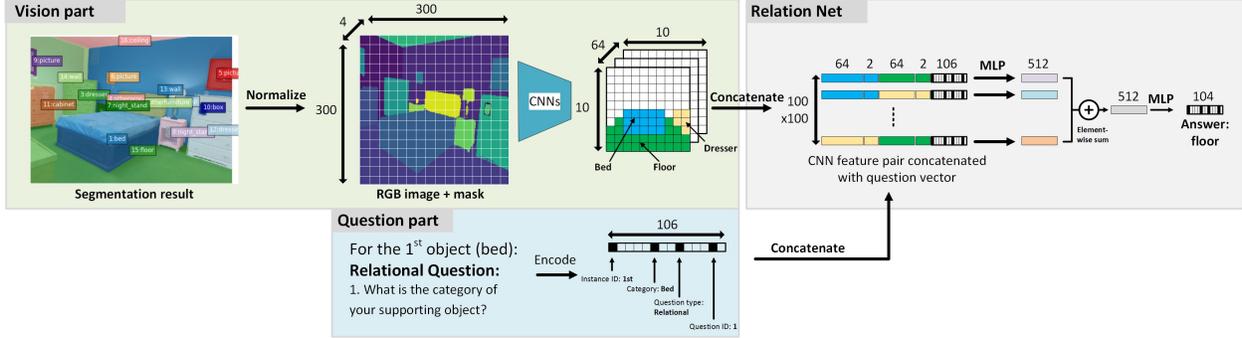
Figure 4: Relation Network for support inference. The whole architecture consists of three parts. The vision part and the question part are responsible for encoding object images and related questions separately, and the Relation Network answers these questions based on the image features.

segmentation results, which instance is supporting object A? Is it supported from below or behind? With this insight, we configure a Relation Network to answer these support relationship questions by linking image features. Our network is designed as shown in Figure 4. The upstream of the Relation Network consists of two parts which encode visual images (with masks) and questions respectively.

In the **Vision** part, the RGB image (color intensities, 3-channel) is normalized to $[0, 1]$ and appended with its mask (instance labels, 1-channel), followed by a scale operation to a 300x300x4 matrix. We then generate 10x10x64 CNN feature vectors after convolutional operations. In the **Question** part, for each object instance, we customize our relational reasoning by answering two groups of questions: non-relational and relational; four questions for each. Taking the bed in Figure 4 as an example, the related questions and corresponding answers are encoded as shown in Figure 5. We design the four relational questions for support inference, and the other four non-relational questions as regularization terms to make our network able to identify the target object we are querying. In our implementation, we train the network on NYU v2 [21]. In a single image, maximal 60 indoor instances with 40 categories are considered. Therefore, for the $i$-th object which belongs to the $j$-th category, we encode the $k$-th question from the $m$-th group to a 106-d (106=60+40+4+2) binary vector.

The outputs of the **Vision** and the **Question** parts are concatenated. We represent the 10x10x64 CNN features by 100 of 64-d feature vectors, and form all possible pairs of these feature vectors into 100x100 pairs. The 100x100 feature pairs are appended with their 2D coordinates (2-d) and exhaustively concatenated with the encoded question vector (106-d), then go through two multi-layer perceptrons to answer the questions (see network specifications in Appendix A.3). For each question, the Relation Network outputs a scalar between 0 and 103. We decode it into a human-language answer by indexing the lookup table as illustrated in Figure 5. The correct rate on the testing dataset of NYU v2 reaches 80.62% and 66.82 % on non-relational and relational questions respectively.

In our experiment, we observe that the numbering of instance masks is randomly given from the object segmentation, which undermines the network performance on the first relational question (see Figure 5). In our implementation, we use the last three relational questions to predict the category of the supporting object and the support type, and keep the first one as a regularization term. The exact supporting instance can be identified by maximizing the prior supporting probability between the target object and its neighbors:

$$O_{j^*} = \operatorname*{argmax}_{O_j \in \mathcal{N}(O_i)} P(\mathcal{C}(O_j)|\mathcal{C}(O_i), T_k), \mathcal{C}(O_j) \in \mathcal{SC}(O_i), \tag{1}$$

where $O_i$ and $\mathcal{N}(O_i)$ respectively denote the $i$-th object and its neighboring instances (layout instances are neighbors to all objects). $\mathcal{C}(O_j)$ represents the category label of object $O_j$. $\mathcal{SC}(O_i)$ indicates the top-5 (in our experiment) category candidates of $O_i$'s supporting object, and $T_k$ denotes the support type, $k = 1, 2$. Hence $P(\mathcal{C}(O_j)|\mathcal{C}(O_i), T_k)$ means the probability of $\mathcal{C}(O_j)$ supporting $\mathcal{C}(O_i)$ by $T_k$. The prior probability P is obtained by counting from the training data (see Appendix B for details). The supporting instance is represented by $O_{j^*}$. This process can improve the testing accuracy on the four relational questions by a large margin (from 66.82% to 82.74%).
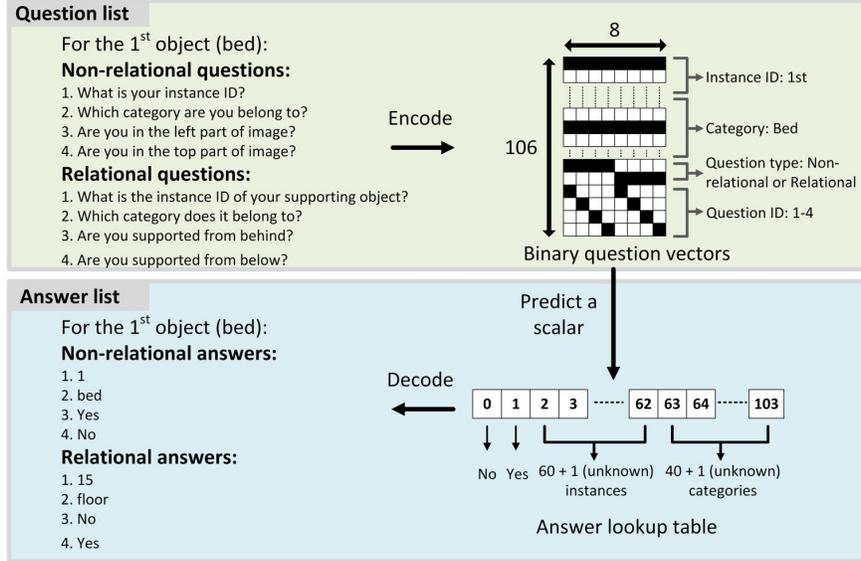
Figure 5: Questions and answers for training

## 6 Global Scene Optimization

The final process is composed of two steps: scene initialization and contextual refinement. The first step initializes camera, 3D layout and object properties. The second step involves an iterative refinement to pick correct object CAD models and fine-tune their sizes, locations and orientations with support relation constraints.

### 6.1 Scene Initialization

**Camera-layout Joint Estimation**    The camera-layout estimation is illustrated in Figure 6. We jointly estimate camera parameters and a refined room layout by minimizing the angle deviations between the layout lines and vanishing lines in images (see Part I in Figure 6). We firstly detect line segments from both the original image and the layout label map using LSD [41] and support vector machine (SVM) respectively. With the initialized camera parameters, orthogonal vanishing points are detected with the strategy proposed by [42]. The quality of vanishing points is scored by the count and length of the line segments they cross through. Longer line segments (like layout lines) would contribute more and guide the orthogonal vanishing lines in alignment with room orientation (see Part I in Figure 6). However, an improper camera initialization, particularly in cluttered environments, would often cause faulty estimation of 3D room layout [7]. We include iterations to improve the camera parameters from the detected line segments and produce a refined room layout simultaneously.
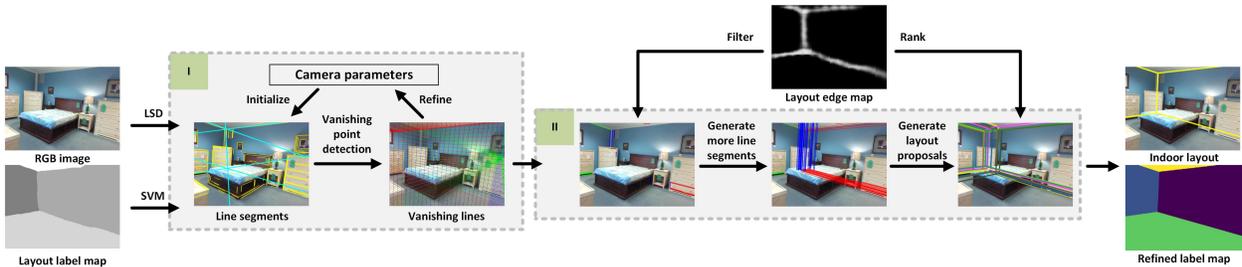


Figure 6: Camera-layout joint estimation. The camera parameters and vanishing points are jointly optimized in Part I, which leads to generate room layout proposals in Part II. The optimal layout is decided by the maximal probability score in layout edge map.

We denote the three orthogonal vanishing points by $\{\mathbf{vp}_i\}$, and the line segment set that (nearly) crosses through $\mathbf{vp}_i$ as $\mathcal{L}(\mathbf{vp}_i), i = 1, 2, 3$. Both of them are expressed by homogeneous coordinates. Similar to K-Means clustering, for the $i$-th cluster $\mathcal{L}(\mathbf{vp}_i)$, we re-estimate a new vanishing point $\mathbf{vp}_i^*$ by decreasing its distances to line segments in $\mathcal{L}(\mathbf{vp}_i)$.
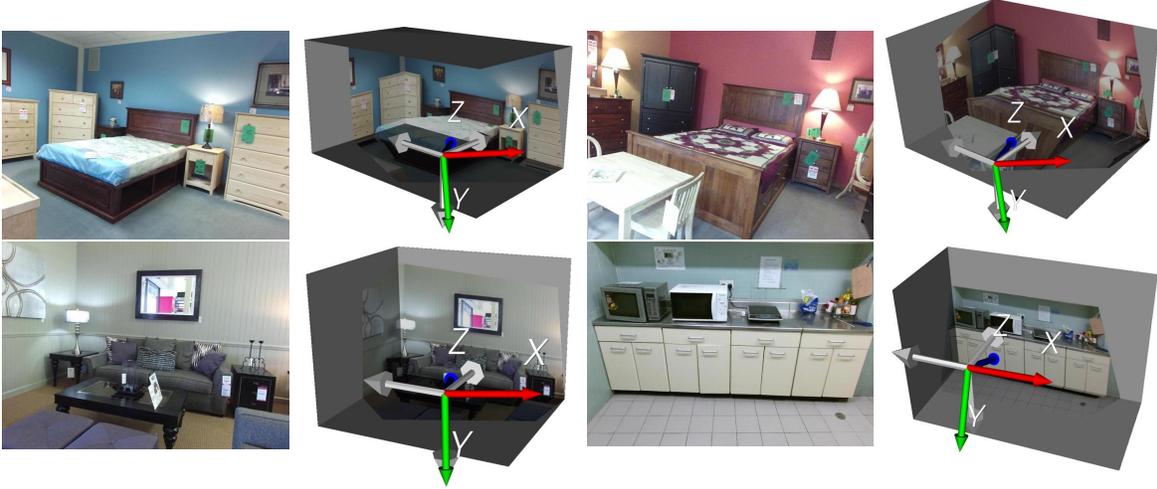
Figure 7: 3D room layout with camera orientation (left: original image, right: 3D layout). The colored arrows represent the camera orientation. The gray arrows respectively point at the floor and walls, which indicates the room layout orientation.

This problem can be formulated as:

$$
\mathbf{vp}_i^* = \underset{\mathbf{vp}_i}{\arg\min}\ \boldsymbol{\varepsilon}^{\mathrm{T}}\boldsymbol{\varepsilon},
$$
$$
[\boldsymbol{l}_1, \boldsymbol{l}_2, ..., \boldsymbol{l}_{N_i}]^{\mathrm{T}}\mathbf{vp}_i = \boldsymbol{\varepsilon}, i = 1, 2, 3, \tag{2}
$$

where $\boldsymbol{l}_k$ denotes the coordinates of a line segment in cluster $\mathcal{L}(\mathbf{vp}_i)$, $k = 1, 2, ..., N_i$. $N_i$ is the capacity of $\mathcal{L}(\mathbf{vp}_i)$. We solve it with the eigen decomposition to obtain the eigen vector corresponding to the smallest eigen value of $[\boldsymbol{l}_1, \boldsymbol{l}_2, ..., \boldsymbol{l}_{N_i}]^{\mathrm{T}}[\boldsymbol{l}_1, \boldsymbol{l}_2, ..., \boldsymbol{l}_{N_i}]$ as the updated $\mathbf{vp}_i$. After that, camera parameters can be updated with the renewed vanishing points by [43]. With this strategy, the vanishing points and camera parameters can be jointly optimized as each of them iteratively converges.

To obtain the optimal indoor layout (see Part II in Figure 6), the line segments that are not located in the layout edge map (high-intensity area) are removed, and we infer more line segments by connecting vanishing points with intersections of line segments from different clusters. More layout proposals can be generated by extensively combining these line segments (see this work [14] for more details). We use the layout edge map to score each pixel in layout proposals and obtain the optimal one with the maximal sum. As the vanishing points provide the room orientation [42], we fit the indoor layout using a 3D cuboid, with the position of a room corner and layout sizes as optimization variables [27]. Then the camera intrinsic and extrinsic parameters can be estimated. Samples of 3D room layout with calibrated cameras are shown in Figure 7.

**Model Initialization** Model retrieval (see Section 4) provides CAD models and orientations for indoor objects. In this part, we introduce single-view geometry combining with support relationship to estimate object sizes and positions with considering object occlusions. The room layout and vanishing points obtained in Section 6.1 are used to measure the height of each object. The whole process is illustrated in Figure 8.

Taking the nightstand and lamp in Figure 8 as examples, the object $\mathrm{O}_i$ (lamp) is supported by $\mathrm{O}_j$ (nightstand) from below. We denote the 2D mask of $\mathrm{O}_j$ by $\mathbf{M}_j$. $\mathbf{vp}_\mathrm{v} \in \mathcal{R}^2$ is the vertical vanishing point on the image plane. For $\mathbf{M}_j$, we get its height line by scanning the mask boundary with rays originated from $\mathbf{vp}_\mathrm{v}$ (see Figure 8(left)). Each ray connects a pixel on the mask boundary with $\mathbf{vp}_\mathrm{v}$. We estimate the Gaussian kernel density of the radian of these rays, and extract the rays whose radian is a local maxima in density. The 'local maximal' ray that holds the longest intersection with the mask boundary is selected, and the longest intersection is taken as the optimal height line of $\mathrm{O}_j$.

To estimate the real height of objects, we introduce single-view geometry for height measurement (see Figure 8(right)). Specifically, we take the room height line as the reference, and map object's height line onto the reference through the vanishing lines. For $\mathrm{O}_i$ (lamp), we denote its top and bottom of the mapped height line by $t_i$ and $b_i$ respectively. $t_r$ and
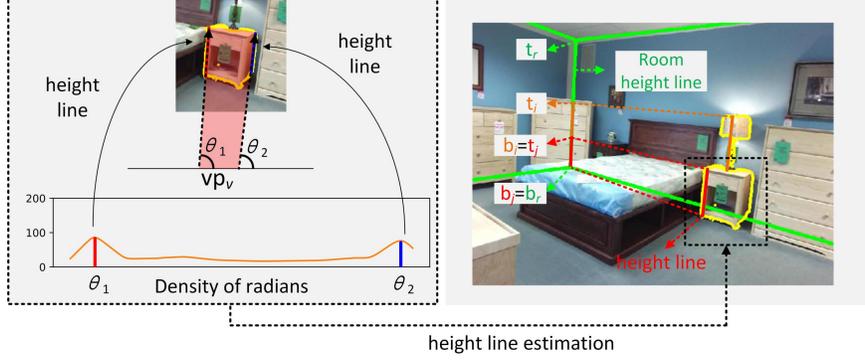
Figure 8: Single-view geometry for object height estimation

$b_r$ respectively indicate the top and the bottom of the room height line. The height of $O_i$ can be calculated by the cross ratio [44]:

$$
\begin{aligned}
H_i &= A_i - A_j, \\
\frac{A_i}{H_r} &= \frac{\|t_i - b_r\|}{\|t_r - b_r\|} \cdot \frac{\|\mathbf{vp_v} - t_r\|}{\|\mathbf{vp_v} - t_i\|},
\end{aligned}
\tag{3}
$$

where $A_i$ and $A_j$ respectively denote the top altitude of $O_i$ and $O_j$ (i.e. the real height of $\overrightarrow{t_i b_r}$ and $\overrightarrow{t_j b_r}$). $O_j$ is supporting $O_i$ from below. Thus $H_i$ is the real height of $O_i$. $H_r$ is the real height of the room (i.e. the real height of $\overrightarrow{t_r b_r}$) and $\| * \|$ represents the Euclidean distance. We use this formula to recursively get the real height of $O_i$ from the difference between the top altitude of $O_i$ and its supporting object $O_j$. Rather than to address their real height individually, this recursive strategy asks for solving equations following the supporting order. It brings us benefits to verify the support type and solve occlusion problems. For example, the support type should be 'support from below' if $H_i$ is larger than zero. Moreover, the bottom of an object ($b_i$) is usually invisible when it is occluded or not segmented out. While in practice, $b_i$ is at the same altitude with $t_j$ if $O_j$ is supporting $O_i$ from below. We replace $b_i$ with $t_j$ in calculations to estimate the real height of each object.

Unlike the 'support from below' scenarios where objects are stacked from the floor following the vertical direction, for objects that are supported from behind, the supporting surfaces are not guaranteed with a fixed normal direction. It would be much more complicated to get a closed-form solution. If $O_i$ is supported by walls (like pictures), we can still get an accurate estimate by Equation 3 (i.e. height difference between $\overrightarrow{t_i b_r}$ and $\overrightarrow{b_i b_r}$). While for other cases (e.g. objects are supported by unknown surfaces), we still use this solution to get a rough estimate first. To ensure a reasonable height estimate, we parse the ScanNet [45] to generate a prior height distribution for each object category and replace those unreasonable estimates with the statistically average (see Appendix B for details).

So far we have obtained the height estimate of each object and its altitude relative to the floor. With the room geometry and the camera parameters obtained in Section 6.1, the 3D location of objects can be estimated using the perspective relation between object masks and its spatial position, we refer readers to this work [46] for more details.

## 6.2 Contextual Refinement

When a room is full of clutter, there could still exist errors in scene initialization, and the aforementioned processes may not be sufficient to solve the scene modeling toward satisfaction. Therefore, a contextual refinement is adopted to fine-tune the CAD models and orientations from candidates (see Section 4). It refines their initial 3D size and position to make the reconstructed scene consistent in semantic and geometric meaning with the indoor context. We formulate this into an optimization problem:

$$
\begin{aligned}
\max_{\theta_i, \boldsymbol{S}_i, \boldsymbol{O}_i, \boldsymbol{p}_i} & \quad \mathtt{IoU}\{\mathtt{Proj}[\boldsymbol{R}(\theta_i) \cdot \boldsymbol{S}_i \cdot \boldsymbol{O}_i + \boldsymbol{p}_i], \mathbf{M}_i\}, \\
\boldsymbol{R}(\theta_i) &= \begin{bmatrix} \cos(\theta_i) & -\sin(\theta_i) & 0 \\ \sin(\theta_i) & \cos(\theta_i) & 0 \\ 0 & 0 & 1 \end{bmatrix}, \; \boldsymbol{S}_i = \begin{bmatrix} s_{i,1} & 0 & 0 \\ 0 & s_{i,2} & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot s_{i,3}, \\
\boldsymbol{p}_i &= [p_{i,1}, \quad p_{i,2}, \quad p_{i,3}]^{\mathrm{T}}, \; i = 1, 2, ..., N.
\end{aligned}
\tag{4}
$$

$\boldsymbol{O}_i$ indicates 3D points in a model candidate of the $i$-th object. All CAD models are initially aligned and placed at the origin of the room coordinate system with the horizontal plane parallel to the floor. $\boldsymbol{S}_i$ is an anisotropic scaling matrix to control the 3D size of $\boldsymbol{O}_i$. $\boldsymbol{R}(\theta_i)$ and $\boldsymbol{p}_i$ are designed to adjust its orientation and position. $\mathtt{Proj}[*]$ denotes the perspective projection to map coordinates from the room coordinate system to the image plane. $\mathtt{IoU}[*]$ is the Intersection over Union operator. $\mathbf{M}_i$ represents the segmented mask of the $i$-th object. Therefore, the target of our contextual refinement is to decide the CAD models $\{\boldsymbol{O}_i\}$ with orientations $\{\theta_i\}$, and adjust their size $\{\boldsymbol{S}_i\}$ and position $\{\boldsymbol{p}_i\}$ to make the 2D projections of those reconstructed objects approximate to our segmentation results. $i = 1, 2, ..., N$ and $N$ indicates the count of segmented objects. We implement the scene refinement with a recursive strategy following the support relation constraints.

**Support constraints from below**    For $\boldsymbol{O}_i$ that is supported by $\boldsymbol{O}_j$ from below, we ask for the geometric center of $\boldsymbol{O}_i$ falling inside the supporting surface, and the bottom of $\boldsymbol{O}_i$ attached above the surface:

$$[\boldsymbol{R}(\theta_i) \cdot \boldsymbol{S}_i \cdot \boldsymbol{O}_i + \boldsymbol{p}_i]^{\mathrm{c}}_{x,y} \geqslant \min[\boldsymbol{O}_j]_{x,y}, \tag{5a}$$

$$[\boldsymbol{R}(\theta_i) \cdot \boldsymbol{S}_i \cdot \boldsymbol{O}_i + \boldsymbol{p}_i]^{\mathrm{c}}_{x,y} \leqslant \max[\boldsymbol{O}_j]_{x,y}, \tag{5b}$$

$$\min[\boldsymbol{R}(\theta_i) \cdot \boldsymbol{S}_i \cdot \boldsymbol{O}_i + \boldsymbol{p}_i]_{z|x,y} \geqslant \max[\boldsymbol{O}_j]_{z|x,y}, \tag{5c}$$

where $[*]^{\mathrm{c}}_{x,y}$ indicates the horizontal coordinate $(x, y)$ of the geometric center, and $[*]_{z|x,y}$ is the altitude value at $(x, y)$.

**Support constraints from behind**    If $\boldsymbol{O}_i$ is supported by $\boldsymbol{O}_j$ from behind, we let $\boldsymbol{O}_i$ to be attached on a side surface of $\boldsymbol{O}_j$'s bounding box. Thus we do not ask for the orientation of $\boldsymbol{O}_i$ as it is consistent with the supporting surface. Considering there are four rectangular side surfaces, for each one, we build a local coordinate system $(\boldsymbol{o}^k_j, \boldsymbol{e}^{k,1}_j, \boldsymbol{e}^{k,2}_j)$ on a vertex $\boldsymbol{o}^k_j$ and a pair of orthogonal edges $(\boldsymbol{e}^{k,1}_j, \boldsymbol{e}^{k,2}_j)$ on these rectangles. $k \in [1, 2, 3, 4]$ indicates one of the four side surfaces, which is decided by solving the target function (4). Support constraints from behind can be written as:

$$0 \leqslant (\boldsymbol{c}_i - \boldsymbol{o}^k_j)^{\mathrm{T}} \cdot \boldsymbol{e}^{k,m}_j \leqslant \|\boldsymbol{e}^{k,m}_j\|^2, \, m = 1, 2, \tag{6a}$$

$$2(\boldsymbol{c}_i - \boldsymbol{o}^k_j)^{\mathrm{T}} \cdot \boldsymbol{n}^k_j = \mathtt{range}[(\boldsymbol{R}(\theta_i) \cdot \boldsymbol{S}_i \cdot \boldsymbol{O}_i)^{\mathrm{T}} \cdot \boldsymbol{n}^k_j], \tag{6b}$$

where

$$\boldsymbol{c}_i = [\boldsymbol{R}(\theta_i) \cdot \boldsymbol{S}_i \cdot \boldsymbol{O}_i + \boldsymbol{p}_i]^{\mathrm{c}}, \tag{6c}$$

$$\boldsymbol{n}^k_j = \boldsymbol{e}^{k,1}_j \times \boldsymbol{e}^{k,2}_j / \|\boldsymbol{e}^{k,1}_j \times \boldsymbol{e}^{k,2}_j\|. \tag{6d}$$

$\boldsymbol{c}_i$ is the geometric center of the updated $\boldsymbol{O}_i$. $\boldsymbol{n}^k_j$ denotes the surface normal (see (6c) and (6d)). Hence, (6a) shows that the projection of $\boldsymbol{c}_i$ along $\boldsymbol{n}^k_j$ should fall inside the supporting surface. $\mathtt{range}[x]$ means $x_{\max} - x_{\min}$. Therefore, (6b) implies that the distance between $\boldsymbol{c}_i$ and the surface should be a half of the object's size along the direction of $\boldsymbol{n}^k_j$. This is to secure the attachment of $\boldsymbol{O}_i$ onto the supporting surface. The only difference from constraint (5) is that the optimization of object orientation turns to choosing a correct supporting surface.

To solve the target function (4), we adopt the exhaustive grid search to decide the exact $\{\boldsymbol{O}_i\}$ and $\{\theta_i\}$. For each grid, BOBYQA method [47] is used to refine $\{\boldsymbol{S}_i\}$ and $\{\boldsymbol{p}_i\}$. We illustrate the convergence trajectory in Figure 9. The results demonstrate that the real height of every objects can be initially estimated before iterative refinement, even though there are heavy occlusions or objects that are not fully segmented. From the IoU curve, 30 iterations for model fine-tuning are enough to recover a whole scene.

# 7   Experiments and Analysis

We present both qualitative and quantitative evaluation of our method with the NYU v2 [21] and SUN RGB-D dataset [48]. All tests are implemented with Python 3.5 on a desktop PC with one TITAN XP GPU and 8 Intel Xeon E5 CPUs. Parameters and network configurations are detailed in Appendix A.

## 7.1   Performance Analysis

We record the average time consumption of each phase for 654 testing samples of NYU v2 (see Table 1). The time cost in modeling a whole scene is related to its complexity. It is expected that modeling a cluttered room with more items costs more time. Object-specific tasks (segmentation, model retrieval) are processed in parallel. On average, it takes 2-3 minutes to process a indoor room of reasonable complexity containing up to 20 detected objects.
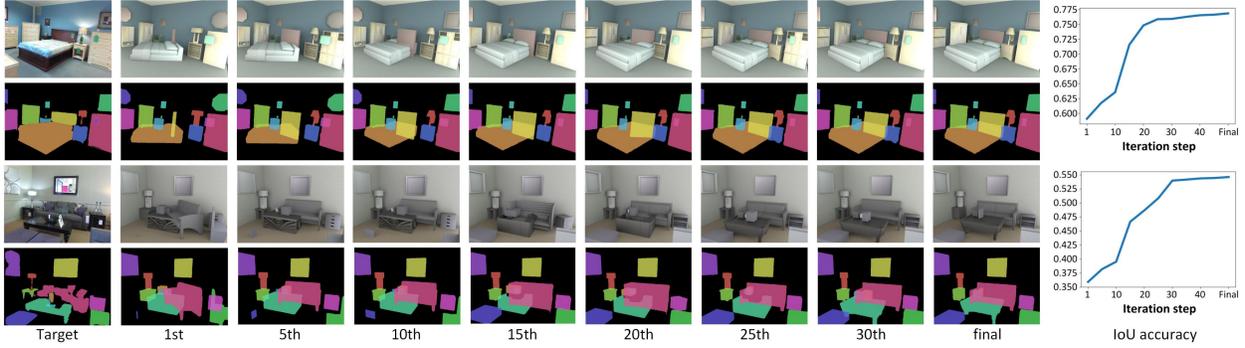
Figure 9: Scene modeling with contextual refinement. The leftmost column presents the original RGB images and the corresponding segmentation. The median part shows the scene modeling results by iterations. The rightmost column illustrates the iteration trajectory of `IoU` values correspondingly.

Table 1: Average time consumption (in seconds) of (1) 2D segmentation + DCRF refining, (2) model retrieval, (3) support inference, (4) camera-layout joint estimation, (5) model initialization and (6) scene modeling. 30 iterations are used in the contextual refinement, and the average number of detected objects is 16 over the 654 testing images.

| Phase | (1) | (2) | (3) | (4) | (5) | (6) | Total |
|---|---|---|---|---|---|---|---|
| Time elapsed | 9.87 | 9.72 | 2.08 | 25.53 | 0.95 | 69.68 | 117.84 |

## 7.2   Qualitative Evaluation

Figure 10 illustrates part of modeling results with different room types and various complexity (randomly picked from the SUN RGB-D dataset, see intermediate results and more samples in Appendix D). The results demonstrate that the detected objects are organized to make the overall presentation consistent with the original images (e.g., object orientation, position and support relationships). The same camera model as the one estimated from each input image is used in rendering, showing both the room layout and camera are reliably recovered with our joint estimation. Benefited from the robust support inference, objects that are heavily occluded or partly visible in the image are predicted with a plausible size.

We compare our outputs with the state-of-the-art works from [9, 7] (see Figure 11). For indoor cases with few objects and occlusions (see Figure 11d, row (1), (2), (4) and (6)), our method extracts more small-size objects (like windows, books, pictures, pillows and lamps) in addition to the main furniture than both methods. This works well with the increasing of scene complexity. Objects that are of low-resolution, hidden or partly out of view can also be captured (see Figure 11a, row (1), (3), (6) and (7)). Both of the two works [9, 7] adopted detection-based methods to locate bounding boxes of objects in a 2D image, which would lose geometric details. Our 'instance segmentation + relational reasoning' approach not only provides more object shape details, but also preserves the relative size between objects. Our context refinement also aligns the recognized models in a meaningful layout driven by the support-guided modeling.

## 7.3   Quantitative Evaluation

We here quantitatively evaluate the 3D room layout prediction, support inference and 3D object placement. Dense modeling of indoor scenes requires the input image to be fully segmented at the instance level. Therefore, we adopt the NYU v2 dataset (795 images for training and 654 images for testing) to assess the tasks of support inference, and use its manually annotated 3D scenes (a subset of the SUN RGB-D annotation dataset) to evaluate the 3D layout prediction and object placement.

**3D Room Layout**   The 3D room layout presents a reference for indoor object alignment and hence influences the object placement. Our method is validated by measuring the average 3D IoU of room bounding boxes between the prediction and the ground-truth [48]. Table 2 illustrates the performance of our method under two configurations: 1. with camera-layout joint estimation and 2. without joint estimation (to estimate camera parameters individually from vanishing points). The results from this ablation experiment show that the strategy of joint estimation consistently outperforms its counterpart in all room types. We also tested the average IoU for 'living rooms' and 'bedrooms' to compare with Izadinia et al. [9]. Our performance reaches 66.08% and Izadinia et al.[9] achieves 62.6% on a subset of SUN RGB-D dataset.

Figure 10: Scene modeling samples on the SUN RGB-D dataset. Each sample consists of an original image (left), the reconstructed scene (raw mesh, middle) and the rendered scene with our estimated camera parameters (right).

Table 2: 3D room layout estimation. Our method is evaluated under two configurations in different room types.

| Room type | bathroom | bedroom | classroom | computer lab | dining room | foyer |
|---|---|---|---|---|---|---|
| IoU (w/o joint) | 30.71 | 39.36 | 47.60 | 20.47 | 46.28 | 54.30 |
| IoU (w/ joint) | **34.90** | **62.86** | **68.23** | **83.21** | **60.41** | **65.59** |
| Room type | kitchen | living room | office | playroom | study room | mean IoU |
| IoU (w/o joint) | 35.37 | 51.34 | 33.49 | 42.91 | 41.93 | 40.10 |
| IoU (w/ joint) | **44.01** | **67.18** | **37.55** | **55.03** | **58.22** | **57.93** |

**Support Inference** The testing dataset from NYU v2 contains 11,677 objects with known supporting instances and support types. Each object is queried with four relational questions. To make fair comparisons with existing methods, we use ground-truth segmentation to evaluate support relations (see [21]). The accuracy of our method is 72.99% at the object level, where a prediction is marked as correct only if all the four questions are correctly answered. This performance reaches the same plateau as existing methods using RGB-D inputs (74.5% by [22] and 72.6% by [21]) and largely outperforms the method using RGB inputs (48.2% by [49]). It demonstrates the feasibility of our Relation Network in parsing support relations from complicated occlusion scenarios without any depth clues.

**3D Object Placement** The accuracy of 3D object placement is tested using manually annotated 3D bounding boxes along with the evaluation benchmark provided by [48], where the mean average precision (mAP) of the 3D IoU between the predicted bounding boxes and the ground-truth is calculated. We align the reconstructed and ground-truth scenes to the same size by unifying the camera altitude, and compare our result with the state-of-the-art [7]. Different from

Figure 11: Comparison with other methods. (a) and (d): The input images. (b) and (e): Reconstructed scenes from other works. The last row is provided by [9], and the remaining results are from [7]. (c) and (f): Our results. All the input images are from the SUN RGB-D dataset.

their work, our method is designed for modeling full scenes with considering all indoor objects, while they adopted a sparsely annotated dataset SUN-RGBD for evaluation with their 30 object categories. As the ground-truth bounding boxes of objects are not fully labeled, we remove those segmented masks that are not annotated to enable comparison under the same configuration. Table 3 shows our average precision scores on the NYU-37 classes [21] (excluding 'wall', 'floor' and 'ceiling'; mAP is calculated with IoU threshold at 0.15). We obtain the mAP score at 11.49. From Huang et al. [7]'s work, they achieved 12.07 on 15 main furniture and 8.06 on all their 30 categories. It shows that our approach achieves better performance in 'smaller' objects, which is in line with the qualitative analysis. The reason could be twofold: 1. a well-trained segmentation network can capture more shape details of objects (e.g. object contour) than using 2D bounding box localization; 2. most human-made objects appear with clear line segments or contours (cabinet, nightstand, dresser, etc.) which benefits our camera-layout joint estimation and model initialization. However, for objects with a rather thin or irregular shape, or under incomplete segmentation (like chair, pillow and lamp et al.), the performance would drop by a small extent.

**Ablation analysis** We implement the ablation analysis to discuss which module in our pipeline contributes most to the final 3D object placement. Two ablated configurations are considered (see Table 3): 1. without camera-layout joint estimation [9], 2. without Relation Network (replaced with prior-based support inference [8]). The mAP scores of the first and second configurations are 5.41 and 8.53 respectively. Our final score is 11.49. It implies that both the

Table 3: 3D object detection. We compare our method under three configurations: 1. without camera-layout joint estimation (w/o joint); 2. without Relation Network (w/o RN); 3. with joint estimation and Relation Network (all). The values show the average precision score on our shared object classes. The column 'others' contains the remaining NYU v2 categories (mAP is averaged by 34 categories, i.e. NYU-37 classes excluding 'wall', 'floor', 'ceiling').

| Method | bathtub | bed | bookshelf | cabinet | chair | desk | door | dresser | fridge | lamp |
|---|---|---|---|---|---|---|---|---|---|---|
| Huang et al. [7] | 2.84 | **58.29** | 7.04 | 0.48 | **13.56** | 4.79 | **1.56** | 13.71 | 15.18 | 2.41 |
| Ours (w/o joint) | 30.83 | 22.62 | 5.83 | 1.82 | 1.12 | 4.31 | 0.68 | 28.53 | 25.25 | 3.12 |
| Ours (w/o RN) | 40.00 | 54.21 | 6.67 | 3.59 | 2.13 | 7.61 | 0.16 | 31.74 | 45.37 | 2.78 |
| Ours (all) | **44.88** | 55.53 | **9.41** | **4.58** | 6.49 | **7.69** | 0.18 | **37.76** | **52.08** | **3.65** |

| Method | nightstand | person | pillow | shelves | sink | sofa | table | toilet | others | mAP |
|---|---|---|---|---|---|---|---|---|---|---|
| Huang et al. [7] | 8.80 | 4.04 | - | - | 2.18 | 28.37 | 12.12 | 16.50 | - | - |
| Ours (w/o joint) | 8.35 | 5.00 | 0.58 | 1.02 | 0.00 | 24.26 | 7.65 | 13.13 | 0.00 | 5.41 |
| Ours (w/o RN) | 32.51 | 8.08 | 0.20 | 3.57 | 0.25 | 31.93 | 7.56 | 10.74 | 0.00 | 8.53 |
| Ours (all) | **32.52** | **18.52** | **1.19** | **33.31** | **3.85** | **33.49** | **13.68** | **31.77** | 0.00 | **11.49** |


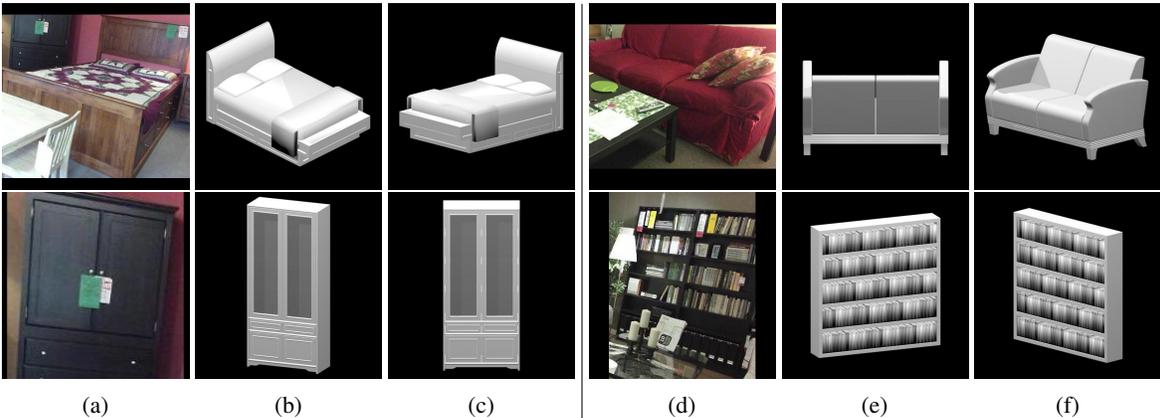
(a)  (b)  (c)  (d)  (e)  (f)

Figure 12: Orientation correction. (a) and (d): The object images. (b) and (e): Matched models from MVRN. (c) and (f): Corrected orientations.

camera-layout joint estimation and relational reasoning contribute to the final performance, and room layout has a higher impact to the object placement in single-view modeling. It is expected that, the orientation and placement of the room layout largely influence the object placement. We also observe that prior-based support inference is more sensitive to occlusions and segmentation quality [8, 22]. When indoor scenes are cluttered, occlusions generally make the supporting surfaces invisible and the segmentation under quality. Unlike the Relation Network, the prior-based method does not take spatial relationship into account and chooses a supporting instance only considering the prior probability, making it more error-prone to complicated scenes.

## 7.4 Discussions

**Improving the Estimation of Object Orientation**    Although the view-based model matching provides an initial guess of object orientation (see Section 4), those deep features are in some cases too abstract to decide sufficiently accurate orientation for trustworthy model initialization. For each object mask, we specifically append a ResNet-34 to predict the orientation angle relative to the camera. It is trained on our dataset considering eight uniformly sampled orientations (i.e. $\pi/4, \pi/2, ..., 2\pi$). However, there is a gap between the renderings (which we used for training) and the real-world images. Rather than conducting full-layer training, we fix the shallowest three layers with the weights pretrained on ImageNet to make our network sensitive to real images. The training data is augmented with coarse drop-out to mimic occlusion effects, and random perspective & affine transformations to mimic different camera poses. The top-1 precision on our testing dataset reaches 91.81% (22342 models for training, 2482 models for testing). Figure 12 illustrates samples from the testing dataset and their predicted orientations. In practice, orientation of some specific models is ambiguous (e.g. symmetric shapes). Top-3 orientation candidates are selected and transformed into the room coordinate system for global scene modeling.
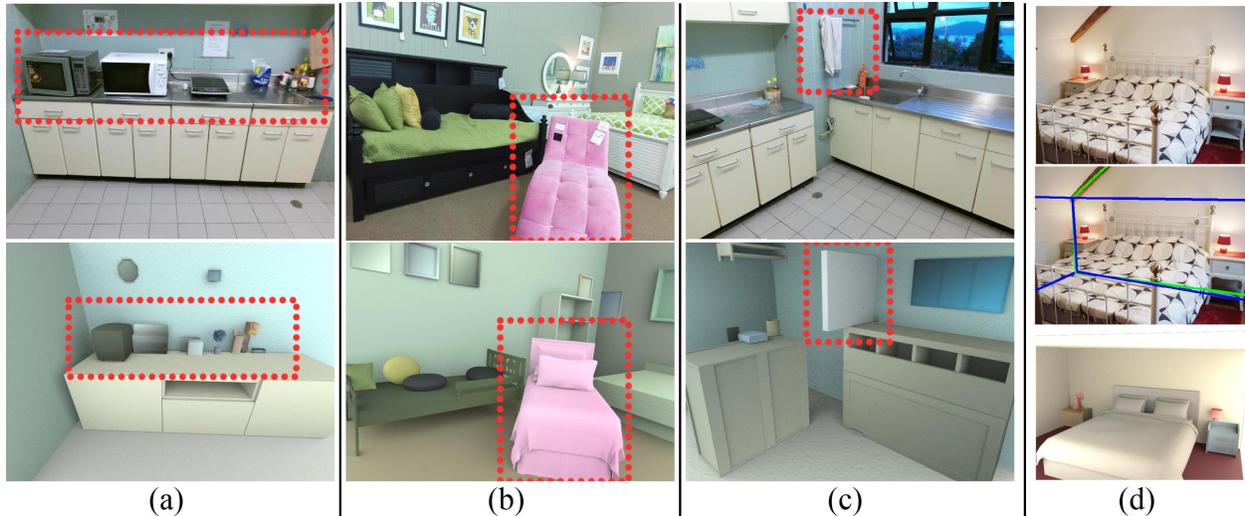
|  |  |  |  |
|---|---|---|---|
| (a) | (b) | (c) | (d) |

Figure 13: Limitation cases. Objects that are segmented with rather few pixels (a), out of our model repository (b) or from 'other category' (right) may not get a proper geometry estimate. For 'non-Manhattan' room layout (d), we fit it with a cuboid. The green and blue lines in (d) respectively represent the 2D room layout and the projection of the 3D layout.

**Limitations**   Our method faces challenges when objects are segmented out with very few pixels (at the minimum of 24x21) which could be too raw for the MVRN to match their shape details. Our CAD model dataset currently contains 37 common categories of indoor objects. Its capacity is limited relative to the diversity of real-world indoor environments. While for unknown objects (labeled as 'other category'), we currently use a cuboid to approximate their shape. Besides, our current method would fit any room layout with a box, which would fail when handling extremely irregular room shapes. Therefore, those reasons above would undermine the IoU accuracy in our contextual refinement, and we illustrate those cases in Figure 13.

## 8   Conclusions and Future Work

We develop a unified scene modeling approach by fully leveraging convolutional features to reconstruct semantic-enriched indoor scenes from a single RGB image. A shallow-to-deep process parses relational and non-relational context into structured knowledge to guide the scene modeling. The experiments demonstrate the capability of our approach in (1) automatically inferring the support relationship of objects, (2) dense scene modeling to recover 3D indoor geometry, with enriched semantics and trustworthy modeling results. Our quantitative evaluations further demonstrate the functionality and effectiveness of each substep in producing semantically-consistent 3D scenes.

This work aims at 3D scene modeling through fully understanding scene context from images. There are high-level relational semantics among indoor objects that could be incorporated into the modeling-by-understanding approach, like other complex contact relations (e.g. a person sits on a chair and holds a mug). All these mixed semantics would help our system to better understand and represent the scene context in a meaningful way. It suggests our future work to provide an intelligent scene knowledge structure to configure and deploy them towards scene modeling.

## References

[1] Kang Chen, Yu-Kun Lai, and Shi-Min Hu. 3d indoor scene modeling from rgb-d data: a survey. *Computational Visual Media*, 1(4):267–278, 2015.

[2] Shuhui Bu, Pengcheng Han, Zhenbao Liu, and Junwei Han. Scene parsing using inference embedded deep networks. *Pattern Recognition*, 59:188–198, 2016.

[3] Hui Wei and Luping Wang. Understanding of indoor scenes based on projection of spatial rectangles. *Pattern Recognition*, 81:497–514, 2018.

[4] Mingming Liu, Kexin Zhang, Jie Zhu, Jun Wang, Jie Guo, and Yanwen Guo. Data-driven indoor scene modeling from a single color image with iterative object segmentation and model retrieval. *IEEE transactions on visualization and computer graphics*, 2018.

[5] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2018.

[6] Yan Zhang, Zicheng Liu, Zheng Miao, Wentao Wu, Kai Liu, and Zhengxing Sun. Single image-based data-driven indoor scene modeling. *Computers & Graphics*, 53:210–223, 2015.

[7] Siyuan Huang, Siyuan Qi, Yixin Zhu, Yinxue Xiao, Yuanlu Xu, and Song-Chun Zhu. Holistic 3d scene parsing and reconstruction from a single rgb image. In *European Conference on Computer Vision*, pages 194–211. Springer, 2018.

[8] Yinyu Nie, Jian Chang, Ehtzaz Chaudhry, Shihui Guo, Andi Smart, and Jian Jun Zhang. Semantic modeling of indoor scenes with support inference from a single photograph. *Computer Animation and Virtual Worlds*, 29(3-4):e1825, 2018.

[9] Hamid Izadinia, Qi Shan, and Steven M Seitz. Im2cad. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5134–5143, 2017.

[10] Chen-Yu Lee, Vijay Badrinarayanan, Tomasz Malisiewicz, and Andrew Rabinovich. Roomnet: End-to-end room layout estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4865–4874, 2017.

[11] Xiaojuan Cheng, Jiwen Lu, Jianjiang Feng, Bo Yuan, and Jie Zhou. Scene recognition with objectness. *Pattern Recognition*, 74:474–487, 2018.

[12] Juan Wang, Xiaoming Tao, Mai Xu, Yiping Duan, and Jianhua Lu. Hierarchical objectness network for region proposal generation and object detection. *Pattern Recognition*, 83:260–272, 2018.

[13] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.

[14] Yuzhuo Ren, Shangwen Li, Chen Chen, and C-C Jay Kuo. A coarse-to-fine indoor layout estimation (cfile) method. In *Asian Conference on Computer Vision*, pages 36–51. Springer, 2016.

[15] Tianqiang Liu, Siddhartha Chaudhuri, Vladimir G Kim, Qixing Huang, Niloy J Mitra, and Thomas Funkhouser. Creating consistent scene graphs using a probabilistic grammar. *ACM Transactions on Graphics (TOG)*, 33(6):211, 2014.

[16] Zhenyu Zhang, Chunyan Xu, Jian Yang, Ying Tai, and Liang Chen. Deep hierarchical guidance and regularization learning for end-to-end depth estimation. *Pattern Recognition*, 83:430–442, 2018.

[17] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2650–2658, 2015.

[18] Yangyan Li, Hao Su, Charles Ruizhongtai Qi, Noa Fish, Daniel Cohen-Or, and Leonidas J Guibas. Joint embeddings of shapes and images via cnn image purification. *ACM Transactions on Graphics (TOG)*, 34(6):234, 2015.

[19] Jiajun Wu, Chengkai Zhang, Xiuming Zhang, Zhoutong Zhang, William T Freeman, and Joshua B Tenenbaum. Learning shape priors for single-view 3d completion and reconstruction. *arXiv preprint arXiv:1809.05068*, 2018.

[20] Kun Xu, Kang Chen, Hongbo Fu, Wei-Lun Sun, and Shi-Min Hu. Sketch2scene: sketch-based co-retrieval and co-placement of 3d models. *ACM Transactions on Graphics (TOG)*, 32(4):123, 2013.

[21] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *European Conference on Computer Vision*, pages 746–760. Springer, 2012.

[22] Feng Xue, Shan Xu, Chuan He, Meng Wang, and Richang Hong. Towards efficient support relation extraction from rgbd images. *Information Sciences*, 320:320–332, 2015.

[23] Bo Zheng, Yibiao Zhao, Joey Yu, Katsushi Ikeuchi, and Song-Chun Zhu. Scene understanding by reasoning stability and safety. *International Journal of Computer Vision*, 112(2):221–238, 2015.

[24] Z. Jia, A. Gallagher, A. Saxena, and T. Chen. 3d-based reasoning with blocks, support, and stability. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, June 2013.

[25] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.

[26] James M Coughlan and Alan L Yuille. Manhattan world: Compass direction from a single image by bayesian inference. In *Proceedings of the Seventh IEEE International Conference on Computer Vision*, volume 2, pages 941–947. IEEE, 1999.

[27] Varsha Hedau, Derek Hoiem, and David Forsyth. Recovering the spatial layout of cluttered rooms. In *2009 IEEE 12th international conference on computer vision*, pages 1849–1856. IEEE, 2009.

[28] Zhuo Deng and Longin Jan Latecki. Amodal detection of 3d objects: Inferring 3d bounding boxes from 2d ones in rgb-depth images. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, page 2, 2017.

[29] Siyuan Huang, Siyuan Qi, Yinxue Xiao, Yixin Zhu, Ying Nian Wu, and Song-Chun Zhu. Cooperative holistic scene understanding: Unifying 3d object, layout, and camera pose estimation. In *Advances in Neural Information Processing Systems*, pages 206–217, 2018.

[30] Joseph J Lim, Aditya Khosla, and Antonio Torralba. Fpm: Fine pose parts-based model with 3d cad models. In *European conference on computer vision*, pages 478–493. Springer, 2014.

[31] Jiajun Wu, Tianfan Xue, Joseph J Lim, Yuandong Tian, Joshua B Tenenbaum, Antonio Torralba, and William T Freeman. Single image 3d interpreter network. In *European Conference on Computer Vision*, pages 365–382. Springer, 2016.

[32] Mingming Liu, Yanwen Guo, and Jun Wang. Indoor scene modeling from a single image using normal inference and edge features. *The Visual Computer*, 33(10):1227–1240, 2017.

[33] Moos Hueting, Pradyumna Reddy, Ersin Yumer, Vladimir G. Kim, Nathan Carr, and Niloy J. Mitra. Seethrough: Finding objects in heavily occluded indoor scene images. In *Proceedings of International Conference on 3DVision (3DV)*, 2018. selected for oral presentation.

[34] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[35] Waleed Abdulla. Mask r-cnn for object detection and instance segmentation on keras and tensorflow. `https://github.com/matterport/Mask_RCNN`, 2017.

[36] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *Advances in neural information processing systems*, pages 109–117, 2011.

[37] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.

[38] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015.

[39] Yu-Shiang Wong, Hung-Kuo Chu, and Niloy J Mitra. Smartannotator an interactive tool for annotating indoor rgbd images. In *Computer Graphics Forum*, volume 34, pages 447–457. Wiley Online Library, 2015.

[40] Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Tim Lillicrap. A simple neural network module for relational reasoning. In *Advances in neural information processing systems*, pages 4967–4976, 2017.

[41] Rafael Grompone Von Gioi, Jérémie Jakubowicz, Jean-Michel Morel, and Gregory Randall. Lsd: a line segment detector. *Image Processing On Line*, 2:35–55, 2012.

[42] Xiaohu Lu, Jian Yaoy, Haoang Li, and Yahui Liu. 2-line exhaustive searching for real-time vanishing point estimation in manhattan world. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 345–353. IEEE, 2017.

[43] Jana Košecká and Wei Zhang. Video compass. In *European conference on computer vision*, pages 476–490. Springer, 2002.

[44] Antonio Criminisi, Ian Reid, and Andrew Zisserman. Single view metrology. *International Journal of Computer Vision*, 40(2):123–148, 2000.

[45] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2017.

[46] Wongun Choi, Yu-Wei Chao, Caroline Pantofaru, and Silvio Savarese. Indoor scene understanding with geometric and semantic contexts. *International Journal of Computer Vision*, 112(2):204–220, 2015.

[47] Powell MJ. The bobyqa algorithm for bound constrained optimization without derivatives. *Cambridge NA Report NA2009/06, University of Cambridge, Cambridge*, 2009.

[48] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 567–576. IEEE, 2015.

[49] Wei Zhuo, Mathieu Salzmann, Xuming He, and Miaomiao Liu. Indoor scene parsing with instance segmentation, semantic labeling and support relationship inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5429–5437, 2017.

[50] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

[51] Alberto Garcia-Garcia, Sergio Orts-Escolano, Sergiu Oprea, Victor Villena-Martinez, and Jose Garcia-Rodriguez. A review on deep learning techniques applied to semantic segmentation. *arXiv preprint arXiv:1704.06857*, 2017.

[52] Saurabh Gupta, Ross Girshick, Pablo Arbeláez, and Jitendra Malik. Learning rich features from rgb-d images for object detection and segmentation. In *European Conference on Computer Vision*, pages 345–360. Springer, 2014.

[53] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.

[54] Guosheng Lin, Chunhua Shen, Anton Van Den Hengel, and Ian Reid. Exploring context with deep structured models for semantic segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1352–1366, 2018.

## A   Technical illustrations

The network configurations and parameter decisions involved in our scene modeling are detailed in this part.

### A.1   Indoor scene segmentation

As Mask R-CNN [13] is designed for general instance segmentation, to make it robust in learning from a small indoor dataset (795 images in our case), we augment the training data with a horizontal flip, and train the network by stages. Specifically, the whole training is divided into three phases, we firstly train the Region Proposal Network, Feature Pyramid Network and mask prediction layers with other parts frozen (60 epochs with learning rate at 1e-3), and fine-tune the ResNet by freezing the shallowest four layers (120 epochs with learning rate at 1e-3) followed by an all-layer training (160 epochs with learning rate at 1e-4). In the inference phase, the searched region proposals go through Non-Maximal Suppression to remove overlaps and keep objects with higher classification scores.

### A.2   Model Retrieval

To build the CAD model dataset, we collect 26,695 indoor models covering 37 categories from ShapeNet [37] and SUNCG [50], along with a 'cuboid' category for objects that are labeled as 'other' in NYU v2. We align and render each model from 32 viewpoints for appearance-based matching, with two elevation angles (15 and 30 degrees) and 16 uniform azimuth angles. The Multi-View Convolutional Network [38] is customized with 32 parallel branches of ResNet-50 [34] as feature extractors (with the last layer removed). All those ResNets share the same weights. The deep features outputted from those branches are max-pooled and fully connected for recognition. The full network is pretrained on ShapeNet for shape recognition task. In our scene modeling, the major color texture from object masks is mapped to CAD models for rendering 3D scenes.

### A.3   Relation Network

The whole architecture consists of three parts (see Figure 5): the Vision part, the Question part, and the Relation reasoning part. The Vision part is designed to encode the image and its segmentation by a set of abstract CNN features. The Question part is to rephrase each question into an encoded vector to ensure our system able to understand human

language. The Relational reasoning part is responsible to analyze the image features and answer the corresponding questions. In the Vision part, we adopt five layers of convolutional kernels (3x3x64 for each layer with the stride and padding size at 2 and 1 respectively). Each convolution is followed by a ReLU and a Batch Normalization layer. The input end is a 300x300x4 matrix (the resized image appended with its mask), and it outputs a 10x10x64 feature map which can be seen as 10x10 of 64-dimensional feature vectors. In the Relational reasoning part, we get exhaustive pair combinations of those 10x10 feature vectors. Each pair of combination is concatenated with their 2D image coordinates correspondingly and the question vector. Thus the image features and the question vector are concatenated into 100x100 visual question vectors. All those vectors separately go through four fully-connected layers, and it generates 100x100 512-dimensional vectors. We take element-wise summation of them and output a (104 dimensional) answer vector after walking-through three fully-connected layers. All the three fully-connected layers above consist of 512 hidden neurons, and each layer is followed by a ReLU unit except the final prediction layer. The initial learning rate is at 0.001 with the wight decay rate at 0.5 in every 10 steps. 60 epoches in total are used for training.

### A.4  Global scene optimization

In Section 6, we set the room height at three meters, and the height of every objects are calculated relatively. To ensure that each height estimate is in a reasonable interval, we parse the ScanNet dataset [45] to conclude a prior height distribution for each object category (see Figure 15 - 18). Each sample in this normal distribution represents a height ratio of a real scanned object relative to the room. A height estimate is regarded as outliers if it is outside the $3\sigma$ interval, and should be replaced by the mean value.

The object sizes and positions are fine-tuned with our contextual refinement. In the optimization problem (see Equation (4)), there are six continuous variables (in $\boldsymbol{S}_i$ and $\boldsymbol{p}_i$) we can control in the optimization process with BOBYAQ method. The constraints (5) and (6) have guaranteed that all objects are attached on their supporting surface. Practically, we further constrain the boundary of $\boldsymbol{S}_i$ to make its size only adjustable in a given interval. we use $s_{i,3}$ in $\boldsymbol{S}_i$ to control the aspect ratio of a CAD model, and $s_{i,1}$ and $s_{i,2}$ to decide its horizontal ratio relative to its height. For common objects (labeled as known NYU v2 categories), we opt to set $s_{i,3} \in [0.9, 1.1]$, and $s_{i,1}, s_{i,2} \in [0.8, 1.2]$. For other objects (labeled as 'other furniture' or 'other structure'), 3D cuboid is used for model retrieval. In this case, we set the boundary of the horizontal ratio more flexible as $s_{i,1}, s_{i,2} \in [0.1, 10]$.

## B  Priors for support inference and height estimation

We parse the ScanNet [45] dataset to get the priors about support relationships and object heights. It contains 1,513 real scene scans with 37,831 indoor objects, and those objects are categorized by the same label set with our experiments. We estimate the bounding box of each object and get the height distribution as the Figure 15 - 18 shows. Each sample in these distributions is a ratio number of the object height to the room height. If a height estimate is beyond $[\mu - 3\sigma, \mu + 3\sigma]$ ($\mu$ is the mean value and $\sigma$ is the standard deviation of the corresponding distribution), we replace the estimate with $\mu$ to initialize the object height.

Moreover, we extract the point cloud of objects to obtain support relationships within all of the scans and get one-to-one support relationship priors (with the method in [39]) as the Figure 19 shows. Each block in the two matrices denotes the number of cases that an object (in row) is supported by another object (in column) from below (Figure 19(a)) or behind (Figure 19(b)). Floating objects are removed, and each object must be supported by another object. When multiple support relationships exist, only the primary one is kept (see [39]).

## C  2D object segmentation comparisons with existing works

2D segmentation is designed to provide the object locations in the image. Detection loss in 2D images directly results in their 3D counterparts missing in the final CAD scenes. Besides that, whether an object is segmented with a fine-grained mask would also affect the geometry estimation. With this concern, we measure the Pixel Accuracy (PA), Mean Accuracy (MA) and Intersection over Union (IoU) [51] between the predicted and ground-truth masks to assess our performance on 40 categories in NYU v2 dataset. In testing, we select object masks with detection score greater than 0.5 from Mask R-CNN and layout masks from FCN to fully segment images. Table 4 illustrates the comparison with state-of-the-art methods. The results demonstrate that we achieve higher performance in terms of PA and IoU scores. It is worth noting that we are mostly concerned about the IoU score which is the optimization target of our contextual refinement.

The 2D IoU from Mask R-CNN [13] only reaches 41.6% though it have reached the state-of-the-art. The accuracy of 3D object placement (i.e. 3D IoU) generally should be much lower for the depth ambiguity. Its indeed a bottleneck for

Table 4: Semantic segmentation on NYU v2 (40 classes). IoU* score is the metric we are concerned in the step of contextual refinement.

| Method | Data type | PA | MA | IoU* |
|---|---|---|---|---|
| Gupta et al.[52] | RGB-D | 60.3 | - | 28.6 |
| FCN-32s [53] | RGB | 60.0 | 42.2 | 29.2 |
| FCN-HHA [53] | RGB-D | 65.4 | 46.1 | 34.0 |
| Lin et al.[54] | RGB | 70.0 | **53.6** | 40.6 |
| Our work | RGB | **70.3** | 49.0 | **41.6** |

all kinds of single image based scene reconstruction methods [7, 9]. However, different from 2D segmentation, the physical plausibility in 3D scene modeling (i.e. relative orientations, sizes, and support relations between objects) could affect the visual performance greater, comparing with the impact from object placement accuracy (i.e. 3D IoU).

In our work, there basically are two factors we most concern: plausibility and placement accuracy. On this basis, we found that obtaining trustworthy physical constraints shows better plausibility and takes more semantic meanings (e.g. support relations) than only chasing placement accuracy. We present an example in Figure 14. In indoor scenes, there are 40 object categories (NYU-40 [21]). Except big-size categories like beds, sofas, tables, etc., most objects are thin or small and occupy little spatial volume (see the pictures and windows in Figure 14). In our experiment, we observed that the 3D IoUs between them and the ground-truth are close to zero, because of their 'skinny' size making the IoU metric vulnerable to placement disparities. However, they are still reconstructed with plausible visual performance because their orientations, sizes and support relations are reasonable. That means, a small 3D offset from the ground-truth will largely lower the accuracy of 3D IoU, but would not affect the visual plausibility given reasonable physical constraints (support relations, orientations and object sizes).
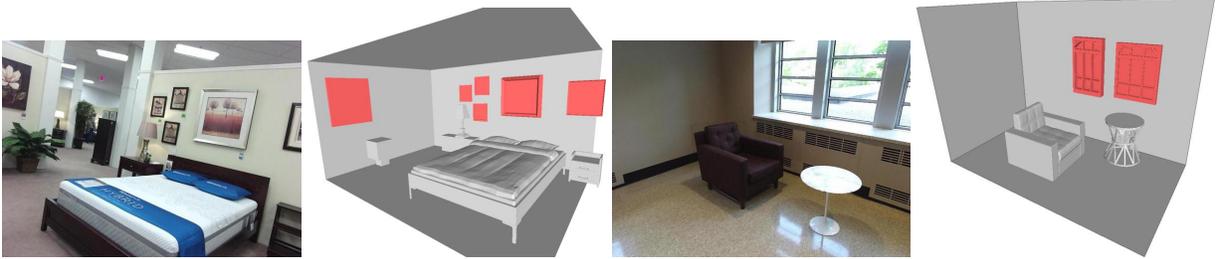


Figure 14: Reconstruction of 'thin' structures.

## D    Intermediate results in scene modeling

We randomly pick around 50 indoor images with different complexity from SUN-RGBD dataset [48]. The modeling results with intermediate outputs are illustrated in Figure 20. The first column shows the input image. The layout edge map and label map are placed in the second and the third column respectively. The fourth column presents the jointly estimated room layout. We illustrate the scene segmentation and the support inference results in the fifth column. Note that the support relationship is represented with an arrow. For example, the red arrow from A to B denotes A supports B from below, and the blue arrow denotes A supports B from behind. We put the modeled scenes in the sixth column (raw scene meshes without texture-mapping and rendering)
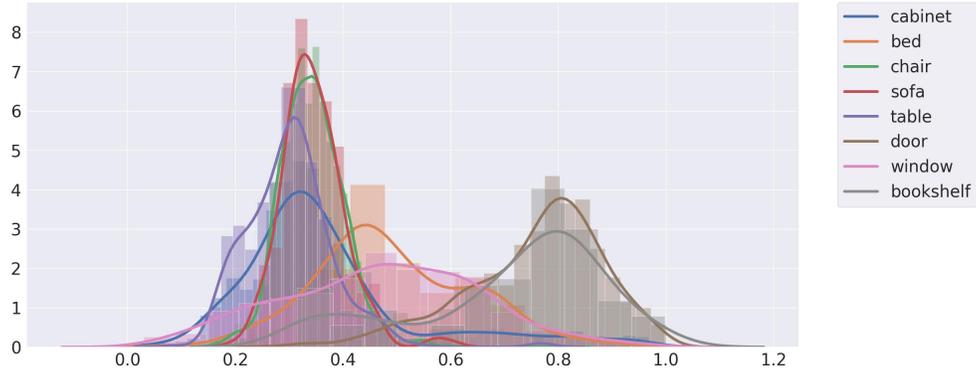
Figure 15: Height distribution for each object category. (1-8 categories)
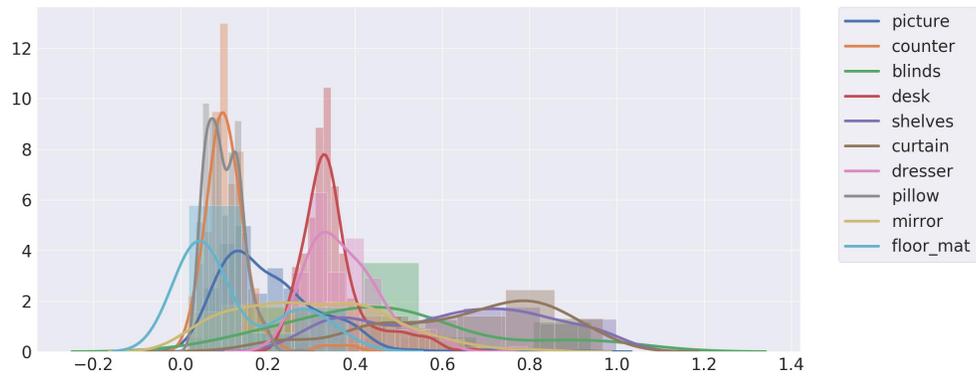


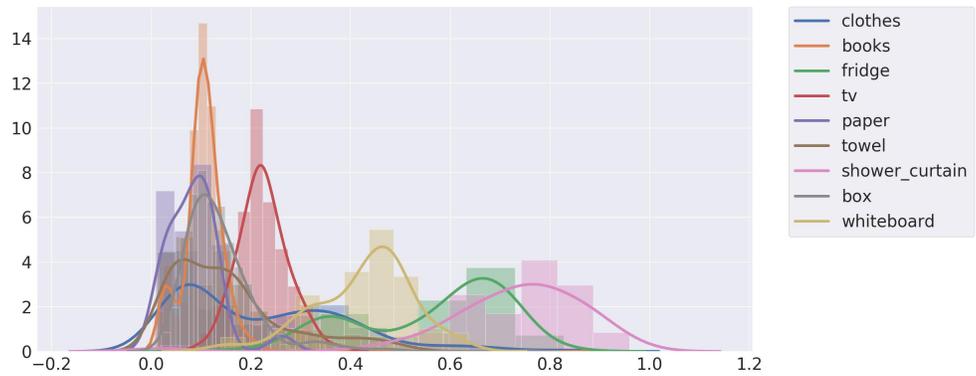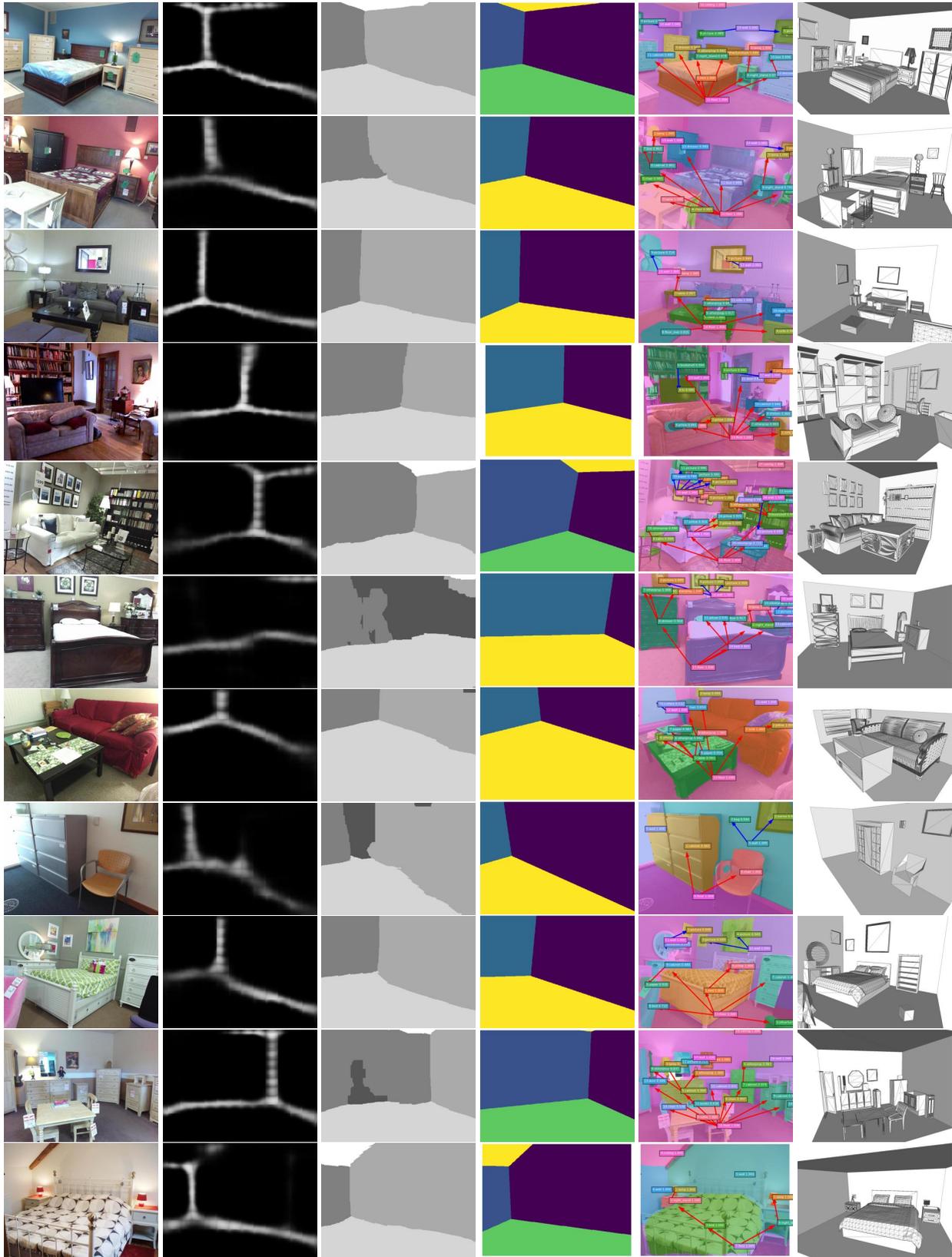Figure 16: Height distribution for each object category. (9-18 categories)



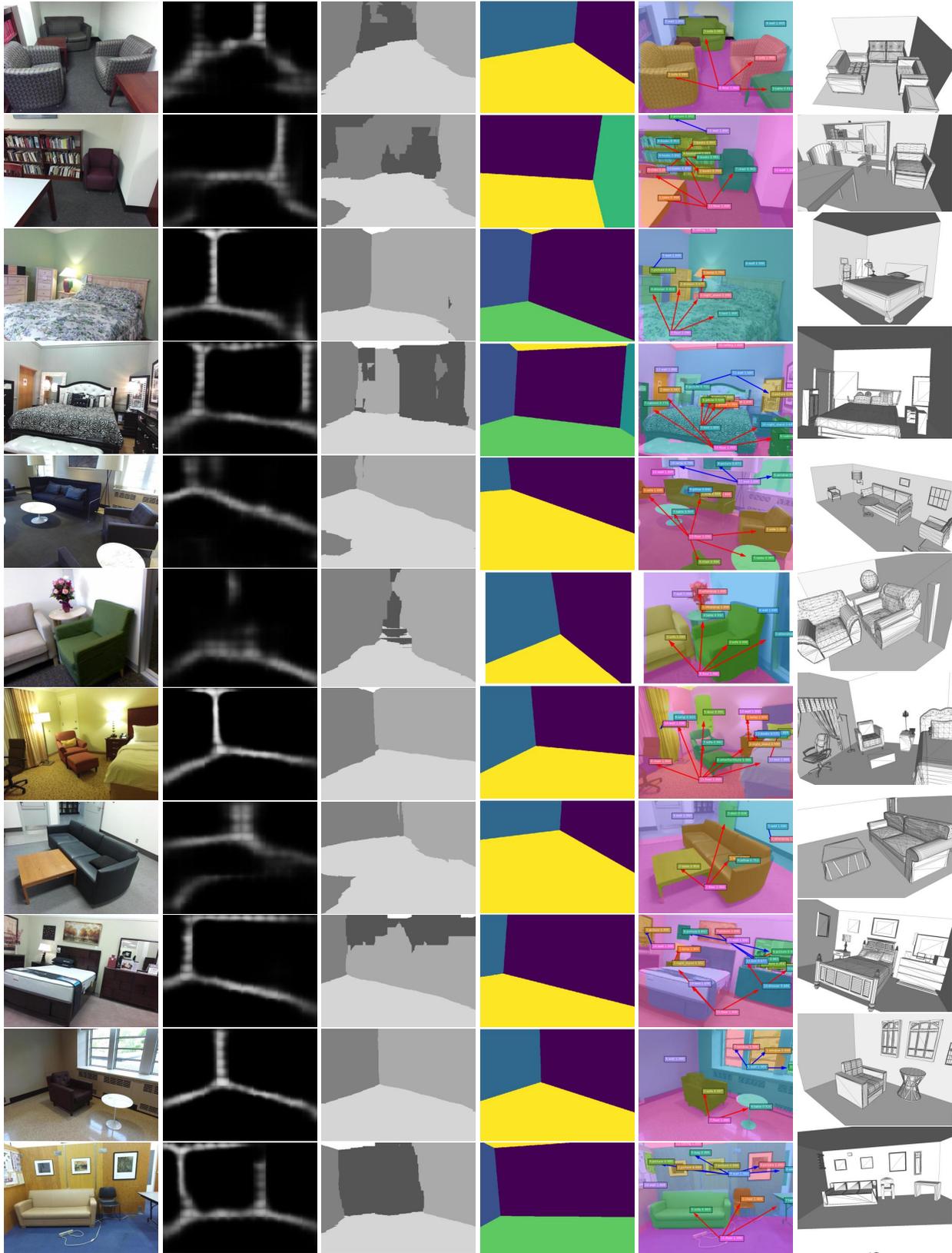Figure 17: Height distribution for each object category. (19-27 categories)

Figure 18: Height distribution for each object category. (28-37 categories)



(a) Support from below



(b) Support from behind

Figure 19: Support relationship priors

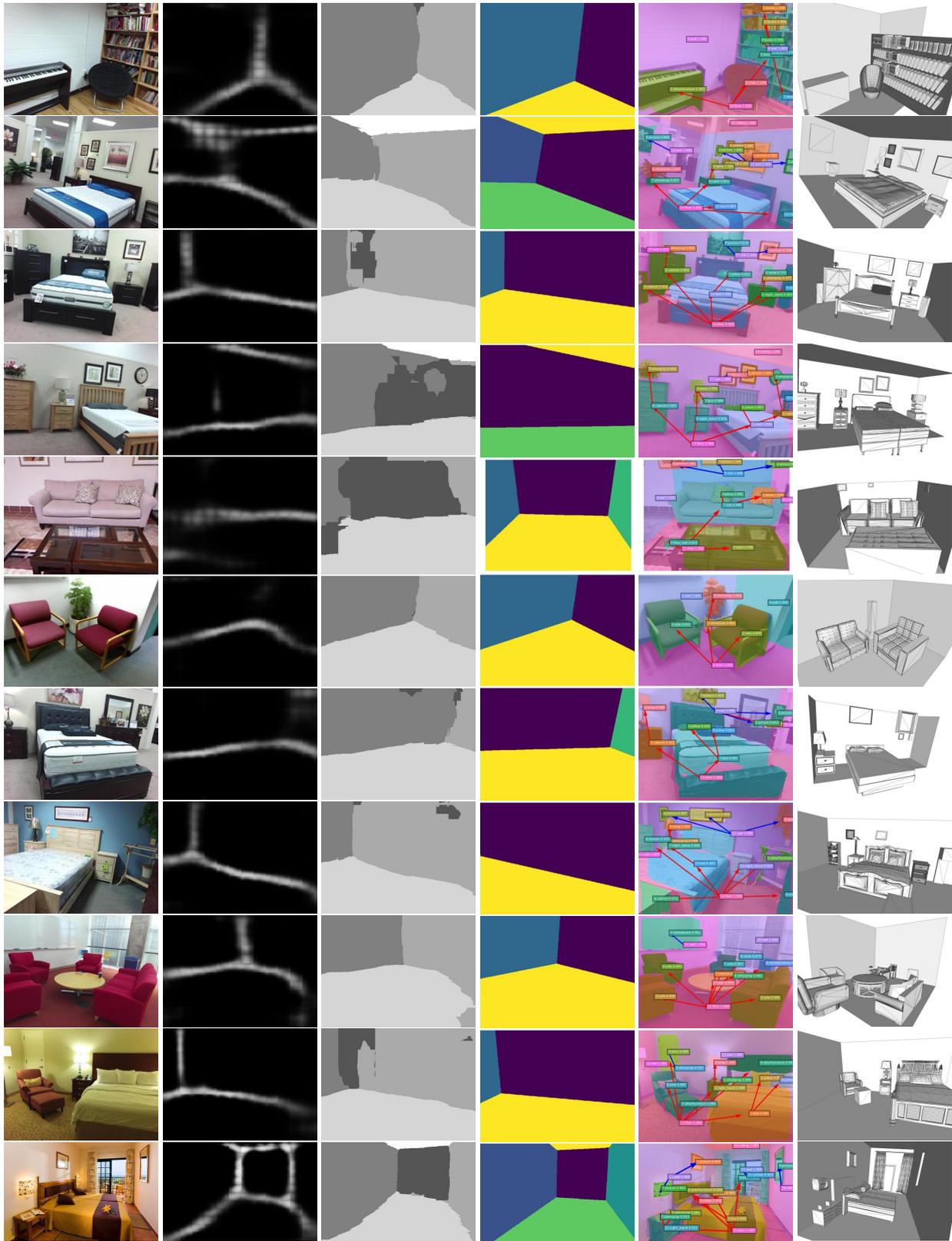(a)　　　(b)　　　(c)　　　(d)　　　(e)　　　(f)

Continue to the next page.

(a)          (b)          (c)          (d)          (e)          (f)

Continue to the next page.
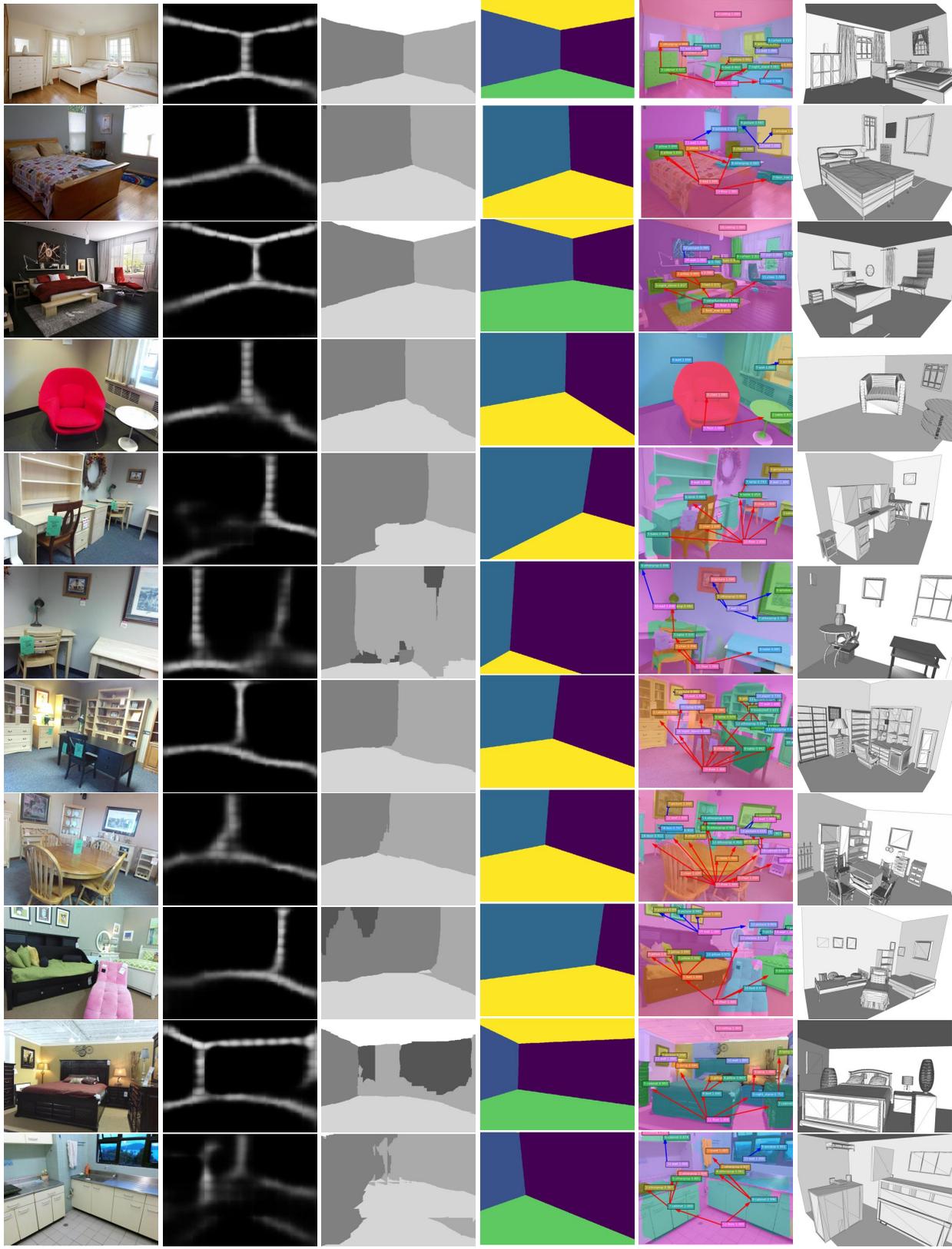
(a)　　　(b)　　　(c)　　　(d)　　　(e)　　　(f)

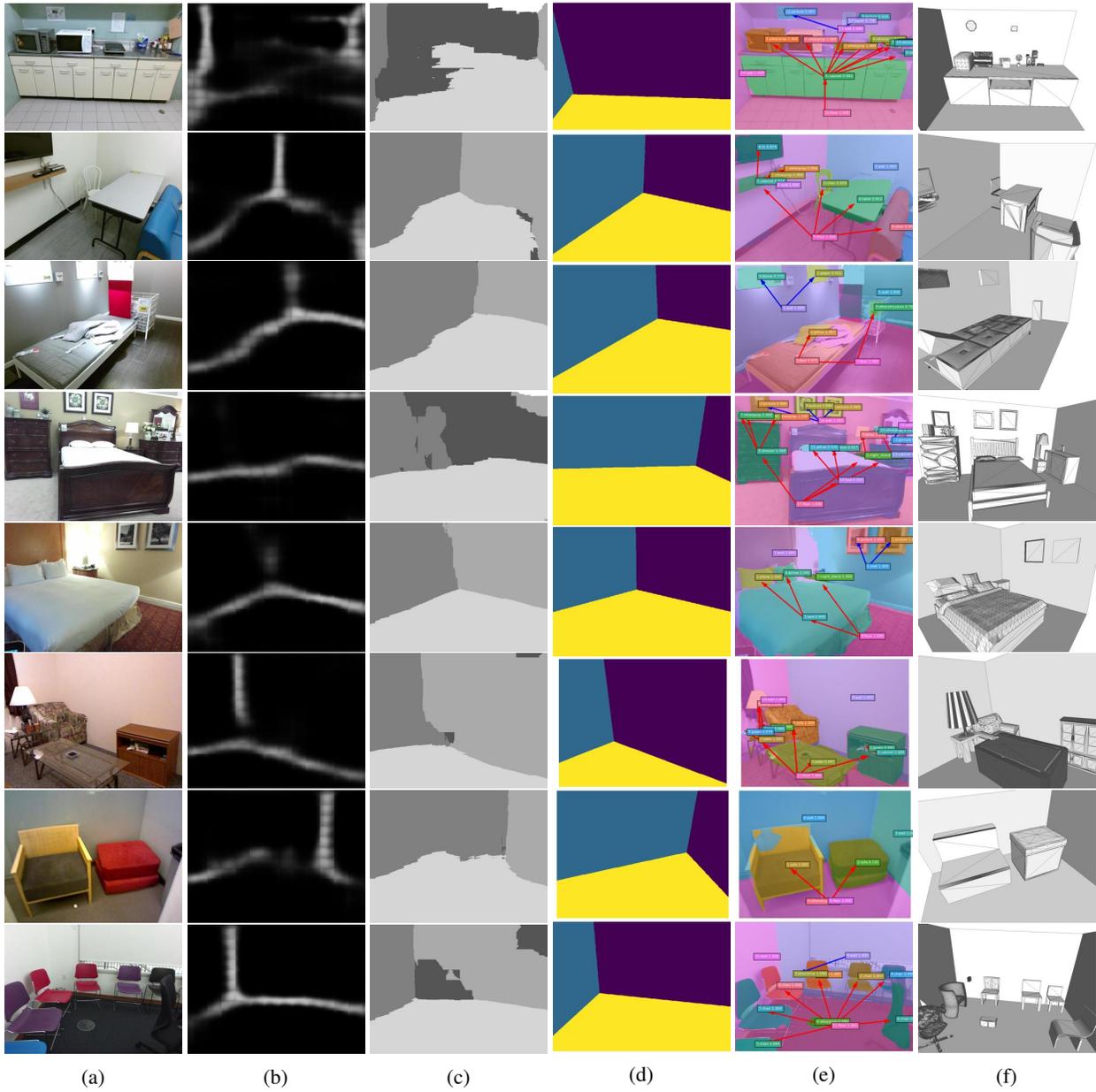Continue to the next page.

(a)   (b)   (c)   (d)   (e)   (f)

Continue to the next page.

Figure 20: Intermediate results in scene modeling.