

Semantic Scene Completion from a Single Depth Image

Shuran Song Fisher Yu Andy Zeng Angel X. Chang Manolis Savva Thomas Funkhouser
Princeton University
<http://sscnet.cs.princeton.edu>

Abstract

This paper focuses on semantic scene completion, a task for producing a complete 3D voxel representation of volumetric occupancy and semantic labels for a scene from a single-view depth map observation. Previous work has considered scene completion and semantic labeling of depth maps separately. However, we observe that these two problems are tightly intertwined. To leverage the coupled nature of these two tasks, we introduce the semantic scene completion network (SSCNet), an end-to-end 3D convolutional network that takes a single depth image as input and simultaneously outputs occupancy and semantic labels for all voxels in the camera view frustum. Our network uses a dilation-based 3D context module to efficiently expand the receptive field and enable 3D context learning. To train our network, we construct SUNCG - a manually created large-scale dataset of synthetic 3D scenes with dense volumetric annotations. Our experiments demonstrate that the joint model outperforms methods addressing each task in isolation and outperforms alternative approaches on the semantic scene completion task. The dataset and code is available at <http://sscnet.cs.princeton.edu>.

1. Introduction

We live in a 3D world where empty and occupied space is determined by the physical presence of objects. To successfully navigate within and interact with the world, we rely on an understanding of both the 3D geometry and the semantics of the environment. Similarly, for a robot, the ability to infer complete 3D shape from partial observations is necessary for low-level tasks such as grasping and obstacle avoidance [33], while the ability to infer the semantic meaning of objects in the scene enables high-level tasks such as retrieval of objects.

With this motivation, our goal is to have a model that predicts both volumetric occupancy (i.e., scene completion) and object category (i.e., scene labeling) from a single depth image of a 3D scene — in this paper we refer to this task as **semantic scene completion** (Figure 1). Prior work is lim-

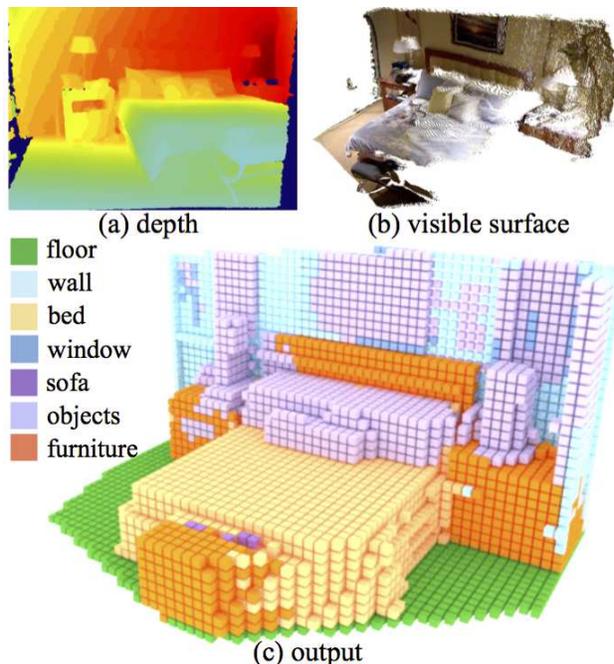


Figure 1. **Semantic scene completion.** (a) Input single-view depth map (b) Visible surface from the depth map; color is for visualization only. (c) Semantic scene completion result: our model jointly predicts volumetric occupancy and object categories for each of the 3D voxels in the view frustum. Note that the entire volume occupied by the bed is predicted to have the bed category.

ited to address only part of this problem as shown in Figure 2: RGB-D segmentation approaches consider only visible surfaces without the full 3D shape [6, 26], while shape completion approaches consider only geometry without semantics [3] or a single object out of context [32, 34].

Our key observation is that the occupancy patterns of the environment and the semantic labels of the objects are tightly intertwined. Therefore, the two problems of predicting voxel occupancy and identifying object semantics are strongly coupled. In other words, knowing the identity of an object helps us predict what areas of the scene it is likely to occupy without direct observation (e.g., seeing the top of a chair behind a table and inferring the presence of a seat and legs). Likewise, having an accurate occupancy pattern for an object helps us recognize its semantic class.

To leverage the coupled nature of the two tasks we jointly train a deep neural network using supervision targeted at both tasks. Given a single-view depth map as input, our semantic scene completion network (SSCNet) produces one of $N+1$ labels for all voxels in the view frustum. Each voxel is labeled as occupied by one of N object categories or free space. Most critically, this prediction extends beyond the projected surface implied by the depth map, thus providing occupancy information for the entire scene.

To achieve this goal there are several issues that must be addressed. First, how do we effectively capture contextual information from 3D volumetric data, where the signal is sparse and lacks high frequency detail? Second, since existing RGB-D datasets only provide annotations on visible surfaces, how do we obtain training data with complete volumetric annotations at scene level?

To address the first issue, we design a 3D dilation-based context module that efficiently expands our network’s receptive field to model the contextual information. We find that a big receptive field is crucial for the task. As demonstrated in Figure 2, looking at the small region of a chair in isolation, it is hard to recognize and complete the chair. However, if we consider the context due to surrounding objects, such as the table and floor, the problem is much easier.

To address the second issue, we construct SUNCG, a large-scale synthetic 3D scene dataset with more than 45622 indoor environments designed by people. All the 3D scenes are composed of individually labeled 3D object meshes, from which we can compute 3D scene volumes with dense object labels through voxelization.

Our experiments with these solutions demonstrate that a method that jointly predicts volumetric occupancy and object semantic can outperform methods addressing each task in isolation. Both the 3D context model learned by our network and the large-scale synthetic training data help to improve performance significantly.

Our main contribution is to formulate an end-to-end 3D ConvNet model (SSCNet) for the joint task of volumetric scene completion and semantic labeling. In support of that goal, we design a dilation-based 3D context module that enables efficient context learning with large receptive fields. To provide the training data for our network, we introduce SUNCG, a manually created large-scale dataset of synthetic 3D scenes with dense occupancy and semantic annotations.

2. Related work

We review related work on RGB-D segmentation, 3D shape completion, and voxel space semantic labeling.

RGB-D semantic segmentation. Many prior works focus on RGB-D image segmentation [6, 26, 29, 15]. However, those methods focus on obtaining semantic labels for only the observed pixels without considering the full shape

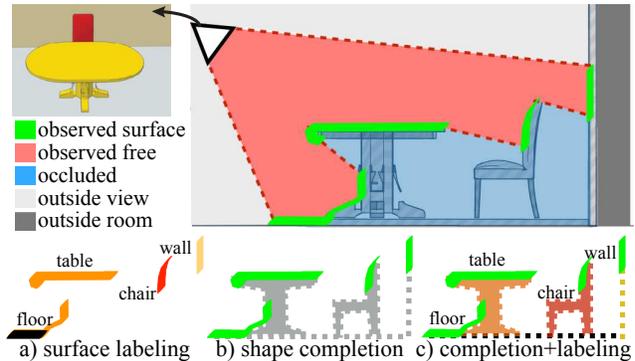


Figure 2. Given a single-view depth observation of a 3D scene the goal of our SSCNet is to predict both occupancy and object category for the voxels on the observed surface and occluded regions.

of the object, and hence cannot directly perform scene completion or predict labels beyond the visible surface.

Shape completion. Other prior works focus on single object shape completion [33, 28, 34, 32]. To apply those methods to scenes, additional segmentation or object masks would be required. For scene completion, when the missing regions are relatively small, methods using plane fitting [22] or object symmetry [13, 19] can be applied to fill in holes. However, these methods heavily rely on the regularity of the geometry and often fail when the missing regions are big. Firman *et al.* [3] show promising completing results on scenes. However, their approach is based purely on geometry without semantics, and thus it produces less accurate results when the scene structure becomes complex.

3D model fitting. One possible approach to obtain the complete geometry and semantic labels for a scene is to retrieve and fit instance-level 3D mesh models to the observed depth map [7, 30, 4, 16, 23, 17, 14]. However, the prediction quality of this type of approach is limited by the quality and variety of 3D models available for retrieval. Naturally, observed objects that cannot be explained by the available models tend to be missed. Or, if the 3D model library is large enough to include all observations, then a difficult retrieval and alignment problem must be solved. Alternatively, it is possible to use 3D primitives such as bounding boxes to approximate the 3D geometry of objects [11, 18, 31]. However, the bounding box approximation limits the geometric detail of the output predictions.

Voxel space reasoning. Another line of work completes and labels 3D scenes, but with separate modules for feature extraction and context modeling. Zheng *et al.* [37] predict the unobserved voxels by physical reasoning. Kim *et al.* [12] train a Voxel-CRF model from labeled floor plans to optimize the semantic labeling and reconstruction for indoor scenes. Hane *et al.* [9] and Blaha *et al.* [1] use joint optimization for multi-view reconstruction and segmentation for outdoor scenes. However, this line of work uses predefined features, and separates the feature learning from the

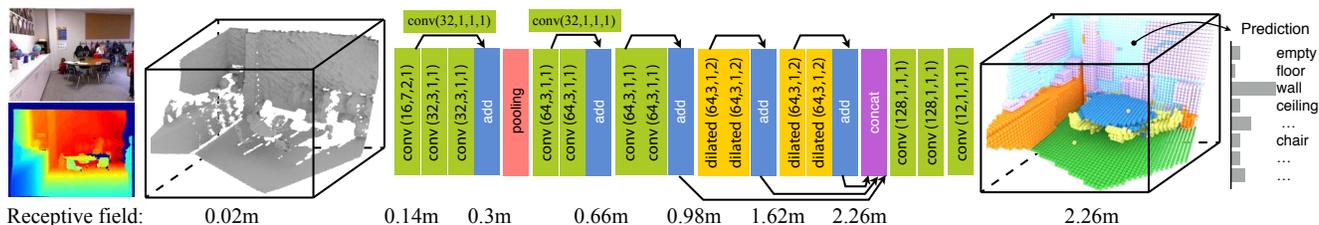


Figure 3. **SSCNet: Semantic scene completion network.** Taking a single depth map as input, the network predicts occupancy and object labels for each voxel in the view frustum. The convolution parameters are shown as (number of filters, kernel size, stride, dilation).

context modeling, and it is expensive for CRF-based models to encode long-range contextual information. In contrast, our model is able to jointly learn the low-level feature representation and high-level contextual information end-to-end from large-scale 3D scene data, directly modeling long-range contextual cues through big receptive field.

Learning from synthetic scene data. Our paper leverages data generated from a large-scale synthetic 3D scene dataset. Although recent works have been focusing on generating segmentation labels for 2D image through rendering synthetic scenes [8, 27], the 3D aspect of such data has not been fully utilized. Existing datasets focus either on objects [2, 34] or a small number of rooms (57 rooms in [8]). In contrast, our dataset is several orders of magnitude larger than existing 3D scene datasets (45,622 houses with 775,574 rooms) providing a diverse set of furniture arrangements manually created by people.

3. Semantic scene completion network

Given a single-view depth map observation of a 3D scene, the goal of our semantic scene completion network is to map the voxels in the view frustum to one of the class labels $C = \{c_0, \dots, c_{N+1}\}$, where N is number of object classes and c_0 represents empty voxels. During training, we render depth maps from virtual viewpoints of our synthetic 3D scenes and voxelize the full 3D scenes with object labels as ground truth. During testing, the observation depth images come from a RGB-D camera.

Figure 3 shows an overview of our processing pipeline. We take a single depth map as input and encode it as a 3D volume. This 3D volume is then fed into a 3D convolutional network, which extracts and aggregates both local geometric and contextual information. The network produces the probability distribution of voxel occupancy and object categories for all voxels inside the camera view frustum.

The following subsections describe the core issues addressed in the design of the system: the data encoding (Section 3.1), network architecture (Section 3.2) and training data generation (Section 4).

3.1. Volumetric data encoding

The first issue we need to address is how to encode the observed depth as input to the network. For the semantic scene completion task, the ideal encoding should directly

represent the 2D observation into the same 3D physical space as the output in a way that is invariant to the view-point projection, and provide a meaningful signal for the network to learn geometry and scene representation. To this end, we adopt Truncated Signed Distance Function (TSDF) to encode the 3D space, where every voxel stores the distance value d to its closest surface, and the sign of the value indicates whether the voxel is in free space or in occluded space. To better suit our task, we make the following modifications to the standard TSDF.

Eliminate view dependency. Most RGB-D reconstruction pipelines speed up the TSDF computation by using the projective TSDF which finds the closest surface points only in the line of sight of the camera [24]. This projective TSDF is fast to compute, but is inherently view-dependent. Instead, we choose to compute the distance to the closest point anywhere on the full observed surface.

Eliminate strong gradients in empty space. Another issue with TSDF is that strong gradients occur in the empty space along the occlusion boundary between $\pm d_{\max}$. It is possible to eliminate this gradient by removing the sign, however, the sign is important for completion task since it indicates the occluded regions of the scene that need to be completed. To solve this problem we flip the TSDF value d as follows: $d_{\text{flipped}} = \text{sign}(d)(d_{\max} - d)$. This flipped TSDF has the strong gradient on surface, providing a more meaningful signal for the network to learn better geometric features. The different encoding is visualized in Figure 5, and Table 3 shows its impact on performance.

3.2. Network architecture

The network architecture of SSCNet is shown in Figure 3. Taking a high-resolution 3D volume as input, the network first uses several 3D convolution layers to learn a local geometry representation. We use convolution layers with stride and pooling layers to reduce the resolution to one fourth of original input. Then, we use a dilation-based 3D context module to capture higher-level inter-object contextual information. After that, the network responses from different scales are concatenated and fed into two more convolution layers to aggregate information from multiple scales. At the end, a voxel-wise softmax layer is used to predict the final voxel label. Several shortcut connections

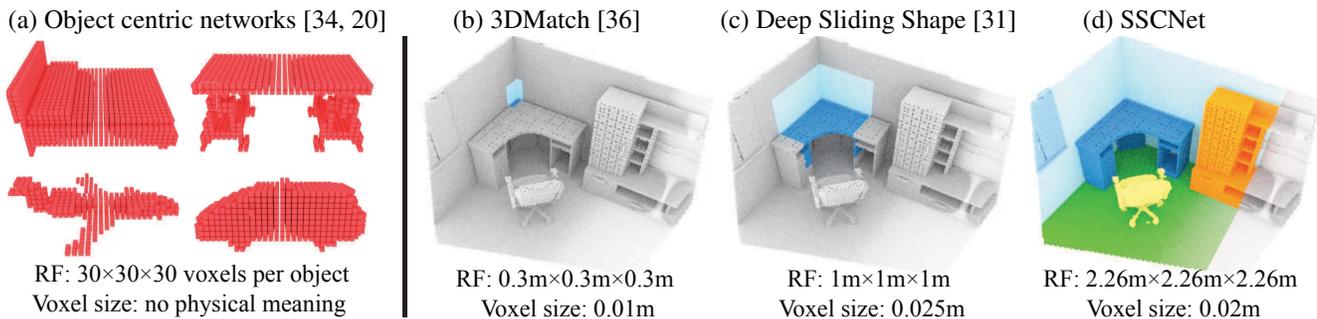


Figure 4. **Comparison of receptive fields and voxel sizes between SSCNet and prior work.** (a) Object centric networks such as [34] and [20] scale objects into the same 3D voxel grid thus discarding physical size information. In (b)-(d), colored regions indicate the effective receptive field of a single neuron in the last layer of each 3D ConvNet. With the help of 3D dilated convolution SSCNet drastically increases its receptive field compared to other 3D ConvNet architectures [31, 36] thus capturing richer 3D contextual information.

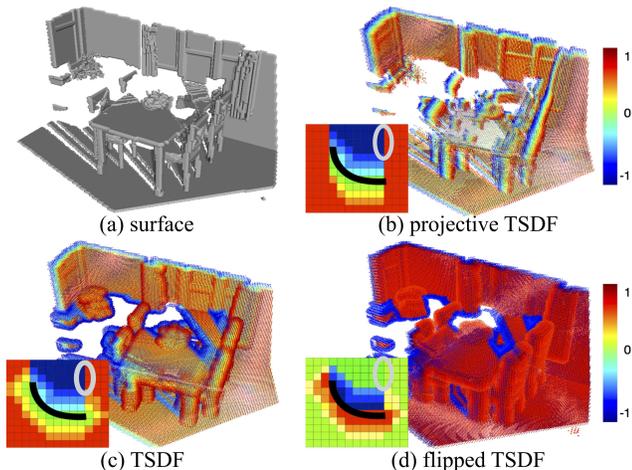


Figure 5. **Different encodings for surface (a).** The projective TSDF (b) is computed with respect to the camera and is therefore view-dependent. The accurate TSDF (c) has less view dependency but exhibits strong gradients in empty space along the occlusion boundary (circled in gray). In contrast, the flipped TSDF (d) has the strongest gradient near the surface.

are added for better gradient propagation. In implementing this architecture, we made the following design decisions:

Input volume generation. Given a 3D scene, we rotate it to align with gravity and room orientation based on Manhattan assumption. The dimensions of the 3D space we consider are 4.8m horizontally, 2.88m vertically, and 4.8m in depth. We encode the 3D scene into a flipped TSDF with grid size 0.02m, truncation value 0.24m, resulting in a $240 \times 144 \times 240$ volume as the network input.

Capturing 3D context with big receptive field. Context can provide valuable information for understanding the scene, as demonstrated by much prior work in image segmentation [35]. In the 3D domain, context is more useful due to a lack of high frequency signals compared to image textures. For example, tabletops, beds, and floors are all geometrically similar to flat horizontal surfaces, so it is hard to distinguish them given only local geometry. However, the relative positions of objects in the scene are a powerful

discriminatory signal. To learn this contextual information, we need to make sure our network has a big enough receptive field. To this end, we extend the dilated convolution presented by Yu and Koltun [35] to 3D. Dilated convolution extends normal convolution by adding a step size when the convolution extracts values from the input before convolving with the kernel. Thus we can exponentially expand the receptive field without a loss of resolution or coverage, while still using the same number of parameters. Figure 4 compares the receptive field size of SSCNet with 3D ConvNet architectures from prior work.

Multi-scale context aggregation. Different object categories have very different physical 3D sizes. This implies that the network will need to capture information at different scales in order to recognize objects reliably. For example, we need more local information to recognize smaller objects like TVs, while we need more global information to recognize bigger objects like beds. In order to aggregate information at different scales we add a layer that concatenates the network responses with different receptive field. We then feed this combined feature map into two $1 \times 1 \times 1$ convolution layers, which allows us to propagate information across responses from different scales.

Data balancing. Due to the sparsity of 3D data, the ratio of empty vs. occupied voxels is around 9:1. To deal with this imbalanced data distribution, we sample the training so that each mini-batch has a balanced set of empty and occupied examples. For each training volume containing N occupied voxels, we randomly sample $2N$ empty voxels from occluded regions for training. Voxels in free space, outside the field of view, or outside the room are ignored.

Loss: voxel-wise softmax. The loss function of the network is the sum of voxel-wise softmax loss $L(p, y) = \sum_{i,j,k} w_{ijk} L_{\text{sm}}(p_{ijk}, y_{ijk})$, where L_{sm} is softmax loss, y_{ijk} is the ground truth label, p_{ijk} is the predicted probability of the voxel at coordinates (i, j, k) over the $N + 1$ classes, where N is the number of object categories and empty voxels are labeled as class 0. The weight w_{ijk} is equal to zero or one based on the sampling algorithm described above.

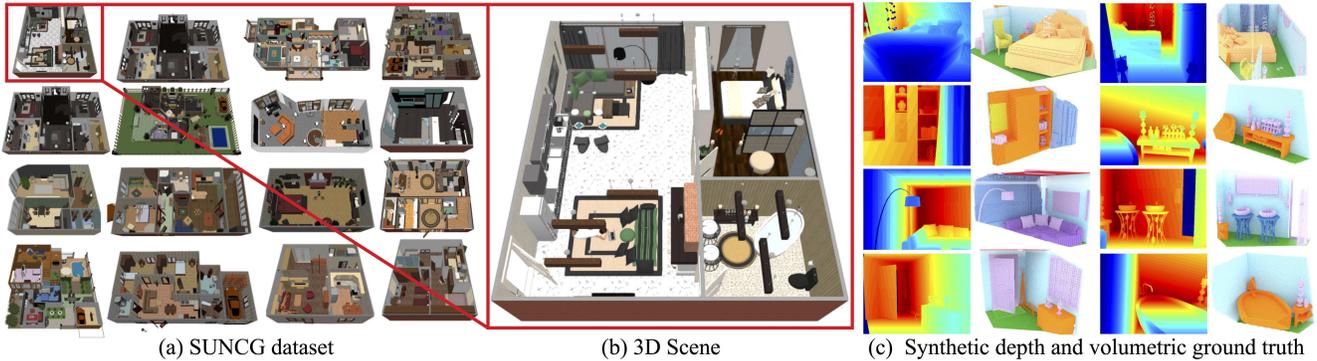


Figure 6. **Synthesizing Training Data.** We collected a large-scale synthetic 3D scene dataset to train our network. For each of the 3D scenes, we select a set of camera positions and generate pairs of rendered depth images and volumetric ground truth as training examples.

Training protocol. We implement our network architecture in Caffe [10]. Pre-training SSCNet on the SUNCG training set takes around a week on a Tesla K40 GPU, and fine-tuning on the NYU dataset takes 30 hours. During training, each mini-batch contains one 3D view volume, requiring 11 GB of GPU memory. To obtain more stable gradient estimates, we accumulate gradients over four iterations and update the weights once afterwards.

4. Synthesizing training data

One of the main obstacles of training deep networks for scene-level dense 3D predictions is the lack of large annotated datasets with dense object semantic annotations at the voxel level. Existing RGB-D datasets with surface reconstructions are subject to occlusions or partial observations, and cannot provide the volumetric occupancy and semantic labels for the entire space at the voxel level. To obtain volumetric occupancy ground truth Firman *et al.* [3] collect a tabletop dataset with reconstructed RGB-D video using KinectFusion [24]. However, this data does not provide semantic labels, and only contains simple tabletop scenarios. In this paper, we present a new large-scale synthetic 3D scene dataset, from which we obtain a large amount of training data with synthetically rendered depth images and volumetric ground truth.

4.1. SUNCG: a large-scale synthetic scene dataset

Our SUNCG dataset contains 45,622 different scenes with realistic room and furniture layouts that are manually created through the Planner5D platform [25]. Planner5D is an online interior design interface that allows users to create multi-floor room layouts, add furniture from a object library, and arrange them in the rooms. After removing duplicated and empty scenes, we ensured the quality of the data with a simple Mechanical Turk cleaning task. During the task, we show a set of top view renderings of each floor and ask turkers to vote whether this is a valid apartment floor. We collect three votes for each floor, and consider a floor valid when it has at least two positive votes. In the end, we have 49,884 valid floors, with contain 404,058 rooms

and 5,697,217 object instances from 2644 unique object meshes covering 84 categories. We manually labeled the all objects in the library to assign category labels. Figure 6 shows example scenes from the resulting SUNCG dataset. More information can be found in the appendix.

4.2. Synthetic depth map generation

To generate synthetic depth maps that mimic a typical image capturing process, we use a set of simple heuristics to pick camera viewpoints. Given a 3D scene, we start with a uniform grid of locations spaced at 1 m intervals on the floor and not occupied by objects. We then choose camera poses based on the distribution of the NYU-Depth v2 dataset.¹ Then, we render the depth map using the intrinsics and resolution of the Kinect. After that we use a set of simple heuristics to exclude bad viewpoints. Specifically, a rendered view is considered valid if it satisfies the following three criteria: a) valid depth area (depth values in range of 1 m to 8 m) larger than 70% of image area, b) there are more than two object categories apart from wall, ceiling, floor, and c) object area apart from wall, ceiling, floor is larger than 30% of image area. To reduce data redundancy, we pick at most five images from each room. In total we generate 130 K valid views for training our SSCNet.

4.3. Volumetric ground truth generation

Since the 3D scenes in the SUNCG dataset consist of a limited number of object instances, we speed up the voxelization process by first voxelizing each individual object in the library and then transforming the labels based on each scene configuration and view point. Specifically, we first voxelize each object to a $128 \times 128 \times 128$ voxel grid. We set the voxel size s so that the largest dimension of the object is a tight fit to the object bounding box. Thus, s varies between objects due to the difference in object dimensions. We use the binvox [21] voxelizer which accounts for both surface and interior voxels by using a space carving approach.

¹The camera height is sampled from a Gaussian distribution with $\mu = 1.5$ m and $\sigma = 0.1$ m. The camera tilt angle is sampled from a Gaussian distribution with $\mu = -10^\circ$ and $\sigma = 5^\circ$.

method (train)	scene completion			semantic scene completion											
	prec.	recall	IoU	ceiling	floor	wall	win.	chair	bed	sofa	table	tv	furn.	objs.	avg.
Lin <i>et al.</i> (NYU) [18]	58.5	49.9	36.4	0	11.7	13.3	14.1	9.4	29	24	6.0	7.0	16.2	1.1	12.0
Geiger and Wang (NYU) [4]	65.7	58	44.4	10.2	62.5	19.1	5.8	8.5	40.6	27.7	7.0	6.0	22.6	5.9	19.6
SSCNet (NYU)	57.0	94.5	55.1	15.1	94.7	24.4	0	12.6	32.1	35	13	7.8	27.1	10.1	24.7
SSCNet (SUNCG)	55.6	91.9	53.2	5.8	81.8	19.6	5.4	12.9	34.4	26	13.6	6.1	9.4	7.4	20.2
SSCNet (SUNCG+NYU)	59.3	92.9	56.6	15.1	94.6	24.7	10.8	17.3	53.2	45.9	15.9	13.9	31.1	12.6	30.5

Table 1. **Semantic scene completion results on the NYU test set with kinect depth map.**

Given a camera view, we define a $240 \times 144 \times 240$ voxel grid in world coordinates, with scene voxel size equals to 2 cm. Then for each object in the scene, we transform the object voxel grid by translating, rotating and scaling by the object’s transformation. We then iterate over each voxel in the scene voxel grid that is inside the transformed object bounding box, and calculate the distance to the nearest neighbor object voxel. If the distance is smaller than the object voxel size s , this scene voxel will be labeled with this object category. Similarly, we label all voxels in the scene that belong to walls, floors, and ceilings by treating them as planes with thickness equal to one scene voxel size. All remaining voxels are marked as empty space, therefore providing a fully labeled voxel grid for the camera view.

5. Evaluation

In this section, we evaluate our proposed methods with a comparison to alternative approaches and an ablation study to better understand the proposed model. We evaluate our algorithm on both real and synthetic datasets.

For the real data, we use the NYU dataset [29], which contains 1449 depth maps captured from Kinect. We obtain the ground truth by voxelizing the 3D mesh annotations from Guo *et al.* [5], mapping object categories based on Handa *et al.* [8]. The annotations consist of 33 object meshes in 7 categories, other categories approximated using 3D boxes or planes. In some cases, the mesh annotation is not perfectly aligned with depth due to labeling error and the limited set of meshes. To deal with this misalignment, Firman *et al.* [3] propose to use rendered depth map from the annotation for testing. However, by rendering the overly simplified meshes, geometric detail is lost especially in cases where objects are represented as a box. Therefore, we test with both rendered depth maps and the originals.

For synthetic data, we created a test set from SUNCG which has objects with detailed geometry, and for which the depth map and ground truth volumes are perfectly aligned. The SUNCG test set are rendered from 184 scenes that are not in the training set.

Evaluation metric. As our evaluation metric, we use the voxel-level intersection over union (IoU) of predicted voxel labels compared to ground truth labels. For the semantic scene completion task, we evaluate the IoU of each object classes on both the observed and occluded voxels. For the

method	training	prec.	recall	IoU
Zheng <i>et al.</i> [37]	NYU	60.1	46.7	34.6
Firman <i>et al.</i> [3]	NYU	66.5	69.7	50.8
SSCNet completion	NYU	66.3	96.9	64.8
SSCNet joint	NYU	75.0	92.3	70.3
SSCNet joint	SUNCG+NYU	75.0	96.0	73.0

Table 2. **Scene completion on the rendered NYU test set as [3]**

scene completion task, we treat all non-empty object class as one category and evaluate IoU of the binary predictions on occluded voxels. Following Firman *et al.* [3], we do not evaluate on voxels outside the view or the room.

5.1. Experimental results

Tables 1 and 2 summarize the quantitative results and Figure 7 shows qualitative comparisons.

Comparison to alternative approaches. In Table 1 we compare on the semantic scene completion task with approaches from Lin *et al.* [18] and Geiger and Wang [4]. Both these algorithms take an RGB-D frame as input and produce object labels in the 3D scene. Lin *et al.* use 3D bounding boxes and planes to approximate all objects. Geiger and Wang retrieve and fit 3D mesh models to the observed depth map at test time. The mesh model library used for retrieval is a superset of the models used for ground truth annotations. Therefore, they can achieve perfect alignments by finding the exact mesh model in a small database. In contrast, our algorithm is based on only depth and does not use additional mesh model at test time. Despite this data disparity, our network produces more accurate voxel-level predictions (30.5% vs. 19.6%). An example of the difference is shown in the third row of Figure 7: both Lin *et al.* and Geiger and Wang’s approaches confuse the sofa as a bed while our network correctly recognizes the sofa. Moreover, since our method does not require the model fitting step it is much faster at 7s compared to 127s per image [4].

Does recognizing objects help scene completion? Previous work has shown scene completion is possible without semantic understanding. We examine to what extent the supervision of object semantics benefits the scene completion task. To do this, we trained a model predicting the occupancy of each voxel by doing binary classification on each

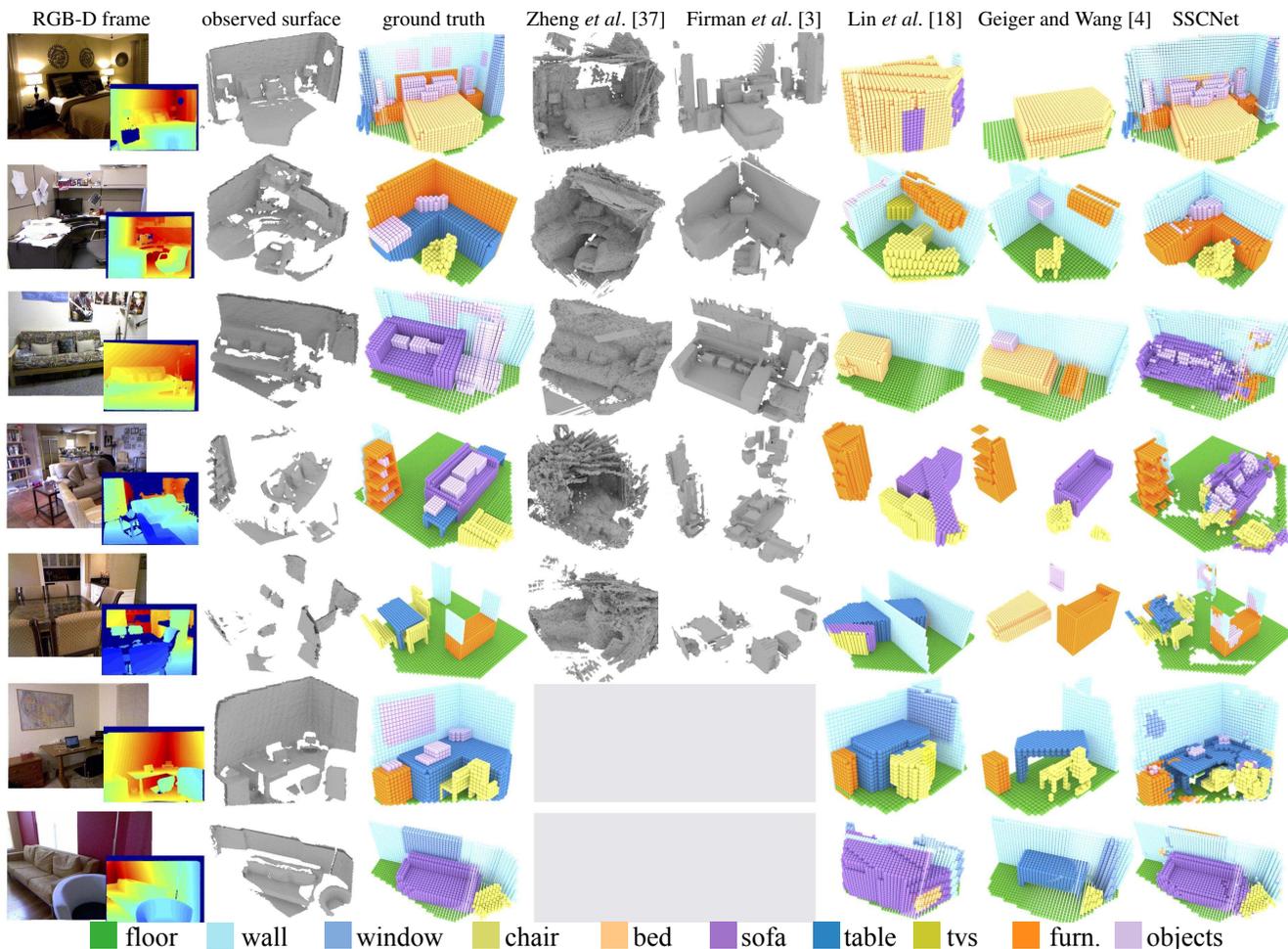


Figure 7. **Qualitative results.** We compare with scene completion results from Zheng *et al.* [37] and Firman *et al.* [3] (on a subset of the test set), and semantic scene completion results from Lin *et al.* [18] and Geiger and Wang [4]. Zheng *et al.* [3] tested on rendered depth, [18] and [4] tested on RGB-D frame from kinect, [3] and SSCNet tested on kinect depth. Overall, SSCNet gives more accurate voxel predictions such as the lamps and pillows in the first row, and the sofa in the third row.

voxel (“empty” vs. “occupied”). We compare the performance of models trained with occupancy and multi-class labeling (see Table 2 [completion] vs. [joint]). We also compare with Firman *et al.* [3] and Zheng *et al.* [37] which both predict binary voxel occupancy based on a single depth map without semantic understanding of the scene. We use the re-implementation of Zheng *et al.*’s approach from Firman *et al.*, which only provides the completion result. We evaluate on the rendered NYU benchmark with the same test images used by Firman *et al.* (randomly picked 200 images from the full test set). While Firman *et al.* produces good results for many cases, their approach fails when the scene becomes complex. For instance, their algorithm fails to complete half of the bed in the first row of Figure 7, and also fails to complete the chairs in the fifth row. In contrast, SSCNet is able to better complete the geometry by leveraging the semantics of the 3D context. This result validates the idea that it is beneficial to understand object semantics in order to achieve better scene completion.

Does scene completion help in recognizing objects? To answer this question, we trained a model with a loss only accounting for semantic labels evaluated on the visible surface and compared with the model trained jointly with labeling and completion (see Table 3 [no completion] vs. [joint]). Even when only evaluating on the visible surface, the model trained with the added supervision of the scene completion task outperforms the model trained only on surface labeling (54.2% vs. 51.2%). This demonstrates that understanding complete object geometry and the 3D context is beneficial for recognizing objects.

Does synthetic data help? To investigate the effect of using synthetic training data, we compared models trained only with NYU and models pre-trained on SUNCG and then fine-tuned on NYU (see Tables 1 and 2 NYU vs. NYU+SUNCG). We see a performance gain by using additional synthetic data especially for the semantic scene completion task having an 10.3% improvement in IoU.

method	encoding	eval	scene completion			semantic scene completion											
			prec.	recall	IoU	ceiling	floor	wall	win.	chair	bed	sofa	table	tv	furn.	objs.	avg.
no completion	flipped	observed	-	-	-	97.2	95.5	61.9	24.6	30.1	55.3	58.9	48.7	14.8	42.1	34.5	51.2
joint	flipped	surface	-	-	-	97.7	94.5	66.4	30.0	36.9	60.2	62.5	56.3	12.1	46.7	33.0	54.2
no balancing	flipped	observed	73.1	95.8	70.8	96.4	85.3	52.1	25.8	16.5	47.1	45.7	28.1	15.3	37.1	19.8	42.7
Basic	flipped	observed	73.4	95.0	70.7	94.6	83.8	47.0	24.0	15.1	38.2	37.2	26.0	0.0	34.8	17.3	38.0
Basic+D	flipped	and	72.2	96.2	70.4	94.7	85.9	47.5	29.2	21.1	50.9	50.7	29.0	21.3	37.2	20.1	44.3
Basic+D+M	proj	occluded	72.0	92.3	67.9	91.6	80.9	45.1	14.6	10.2	39.4	29.8	19.8	0.0	27.4	14.3	33.9
Basic+D+M	tsdf	voxels	74.8	94.0	71.4	95.8	84.4	45.1	17.5	15.2	28.2	37.2	25.6	0.0	28.2	21.9	36.3
Basic+D+M	flipped		76.3	95.2	73.5	96.3	84.9	56.8	28.2	21.3	56.0	52.7	33.7	10.9	44.3	25.4	46.4

Table 3. **Ablation study on SUNCG testset.** First two row shows the evaluation on surface segmentation with and without joint training. The following rows show the evaluation on semantic scene completion task. D: 3D dilated convolution. M: multi-scale aggregation.

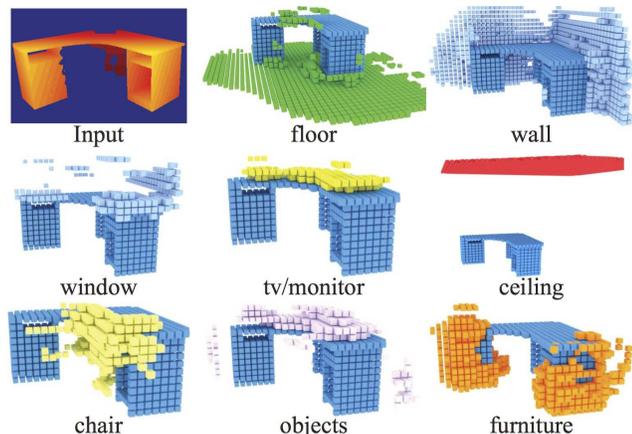


Figure 8. **What 3D context does the network learn?** The first figure shows the input depth map (a desk) and the following figures show the predictions for other objects. Without observing any information for other objects the SSCNet is able to hallucinate their locations based on the observed object and the learned 3D context.

Does a bigger receptive field help? In Table 3, the networks labeled [Basic] and [Basic+D] have the same number of parameter, while in [Basic+D] three convolution layers are replaced by dilated convolution, increasing the receptive field from 1.16 m to 2.26 m. Increasing the receptive field gives the network a opportunity to capture richer contextual information and significantly improve the network performance from 38.0% to 44.3%. To visualize the contextual information learned by the network, we perform the following experiment: given a depth map of a single object we predict labels for all unobserved voxels. Figure 8 shows the input depth and the predictions. Even without observing depth information for other objects SSCNet hallucinates plausible contextual object based on the observed object.

Does multi-scale aggregation help? Comparing the network performance with and without the aggregation layer (see Table 3 [Basic+D] vs. [Basic+D+M]), we observe that the model with aggregation yields higher IoU for both the scene completion and semantic scene completion tasks by 3.1% and 2.1% respectively.

Do different encodings matter? The last three rows in Table 3 compare different volumetric encodings: projective TSDF [proj.], accurate TSDF [tsdf], and flipped TSDF [flipped]. We observe that removing the view dependency by using the accurate TSDF gives a 2.4% improvement in IoU. Making the gradients concentrated on the surface with the flipped TSDF leads to a 10.1% improvement.

Is data balancing necessary? To balance the empty and occupied voxel examples, we proposed to sample the empty voxels during training. In Table 3, [no balancing] shows the performance when we remove the sampling process during training, where we see a drop in IoU from 46.4% to 42.7%.

Limitations. Firstly, we do not use any color information, so objects missing depth such as “windows” are hard to handle. This also leads to confusion between objects with similar geometry or functionality. For example, in the second row of Figure 7 the network predicts the desk as the broader furniture category. Secondly, due to the GPU memory constraints, our network output resolution is lower than that of input volume. This results in less detailed geometry and missing small objects, such as the missed objects on the desk of the second row in Figure 7.

6. Conclusion

In this paper, we introduced SSCNet, a 3D ConvNet for the semantic scene completion task of jointly predicting volumetric occupancy and semantic labels for full 3D scenes. We trained this network on a new large-scale synthetic 3D scene dataset. Experiment results demonstrate that our joint model outperforms methods addressing either component task in isolation, and that by leveraging the 3D contextual information and the synthetic training data, we significantly outperform alternative approaches on the semantic scene completion task.

Acknowledgment This work is supported by Intel, Adobe, and NSF (IIS-1251217 and VEC 1539014/ 1539099). It makes use of data from Planner5D and hardware donated by NVIDIA and Intel. Shuran Song is supported by a Facebook Fellowship. Helpful advice and assistance was provided by Elizabeth R. Bradley, Michael Firman, Yinda Zhang, Matthias Niessner and Angela Dai.

References

- [1] M. Bláha, C. Vogel, A. Richard, J. D. Wegner, T. Pock, and K. Schindler. Large-scale semantic 3d reconstruction: an adaptive multi-resolution model for multi-class volumetric labeling.
- [2] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], 2015.
- [3] M. Firman, O. Mac Aodha, S. Julier, and G. J. Brostow. Structured prediction of unobserved voxels from a single depth image. In *CVPR*, 2016.
- [4] A. Geiger and C. Wang. Joint 3D object and layout inference from a single RGB-D image. In *German Conference on Pattern Recognition (GCPR)*, 2015.
- [5] R. Guo, C. Zou, and D. Hoiem. Predicting complete 3D models of indoor scenes. *ICCV*, 2015.
- [6] S. Gupta, P. Arbelaez, and J. Malik. Perceptual organization and recognition of indoor scenes from RGB-D images. In *CVPR*, 2013.
- [7] S. Gupta, P. A. Arbeláez, R. B. Girshick, and J. Malik. Aligning 3D models to RGB-D images of cluttered scenes. In *CVPR*, 2015.
- [8] A. Handa, V. Patraucean, V. Badrinarayanan, S. Stent, and R. Cipolla. SceneNet: Understanding real world indoor scenes with synthetic data. *CoRR*, abs/1511.07041, 2015.
- [9] C. Hane, C. Zach, A. Cohen, R. Angst, and M. Pollefeys. Joint 3D scene reconstruction and class segmentation. In *CVPR*, pages 97–104. IEEE Computer Society, 2013.
- [10] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [11] H. Jiang and J. Xiao. A linear approach to matching cuboids in RGBD images. In *CVPR*, 2013.
- [12] B. Kim, P. Kohli, and S. Savarese. 3D scene understanding by Voxel-CRF. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1425–1432, 2013.
- [13] Y. M. Kim, N. Mitra, D. Yan, and L. Guibas. Acquisition of 3D indoor environments with variability and repetition. *ACM Trans. Gr*, 2012.
- [14] Y. M. Kim, N. J. Mitra, D. Yan, and L. Guibas. Acquiring 3D indoor environments with variability and repetition. *ACM Trans. Graph. (Proc. SIGGRAPH Asia)*, 31(6), 2012.
- [15] K. Lai, L. Bo, and D. Fox. Unsupervised feature learning for 3D scene labeling. In *ICRA*. IEEE, 2014.
- [16] K. Lai and D. Fox. Object recognition in 3D point clouds using web data and domain adaptation. In *IJRR*, 2010.
- [17] Y. Li, A. Dai, L. Guibas, and M. Nießner. Object detection and classification from large-scale cluttered indoor scans. In *Computer Graphics Forum*, volume 34, 2015.
- [18] D. Lin, S. Fidler, and R. Urtasun. Holistic scene understanding for 3D object detection with RGBD cameras. In *ICCV*, 2013.
- [19] O. Mattausch, D. Panozzo, C. Mura, O. Sorkine-Hornung, and R. Pajarola. Object detection and classification from large-scale cluttered indoor scans. In *Computer Graphics Forum*. Wiley Online Library, 2014.
- [20] D. Maturana and S. Scherer. VoxNet: A 3D convolutional neural network for real-time object recognition. In *IROS*, 2015.
- [21] P. Min. Binvox <http://www.patrickmin.com/binvox/>.
- [22] A. Monszpart, N. Mellado, G. J. Brostow, and N. J. Mitra. Rapter: Rebuilding man-made scenes with regular arrangements of planes. *TOG*, 34(4):103, 2015.
- [23] L. Nan, K. Xie, and A. Sharf. A search-classify approach for cluttered indoor scene understanding. *ACM Trans. Graph. (Proc. SIGGRAPH Asia)*, 31(6), 2012.
- [24] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon. KinectFusion: Real-time dense surface mapping and tracking. In *ISMAR*, 2011.
- [25] Planner5D. <https://planner5d.com/>.
- [26] X. Ren, L. Bo, and D. Fox. RGB-(D) scene labeling: Features and algorithms. In *CVPR*, 2012.
- [27] S. R. Richter, V. Vineet, S. Roth, and V. Koltun. Playing for data: Ground truth from computer games. In *ECCV*, 2016.
- [28] J. Rock, T. Gupta, J. Thorsen, J. Gwak, D. Shin, and D. Hoiem. Completing 3D object shape from one depth image. In *CVPR*, 2015.
- [29] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor segmentation and support inference from RGBD images. In *ECCV*, 2012.
- [30] S. Song and J. Xiao. Sliding Shapes for 3D object detection in depth images. In *ECCV*, 2014.
- [31] S. Song and J. Xiao. Deep sliding shapes for amodal 3D object detection in rgb-d images. In *CVPR*, 2016.
- [32] D. Thanh Nguyen, B.-S. Hua, K. Tran, Q.-H. Pham, and S.-K. Yeung. A field model for repairing 3D shapes. In *CVPR*, June 2016.
- [33] J. Varley, C. DeChant, A. Richardson, A. Nair, J. Ruales, and P. Allen. Shape completion enabled robotic grasping. *arXiv*, 2016.
- [34] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3D ShapeNets: A deep representation for volumetric shapes. In *CVPR*, 2015.
- [35] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016.
- [36] A. Zeng, S. Song, M. Nießner, M. Fisher, and J. Xiao. 3DMatch: Learning the matching of local 3D geometry in range scans. In *arXiv*, 2016.
- [37] B. Zheng, Y. Zhao, J. C. Yu, K. Ikeuchi, and S.-C. Zhu. Beyond point clouds: Scene understanding by reasoning geometry and physics. In *CVPR*, pages 3127–3134. IEEE Computer Society, 2013.