# Mask2CAD: 3D Shape Prediction by Learning to Segment and Retrieve

Weicheng Kuo[1,2], Anelia Angelova[1,2], Tsung-Yi Lin[1,2], and Angela Dai[3]

[1] Google AI
[2] Robotics at Google
[3] Technical University of Munich
{weicheng,anelia,tsungyi}@google.com, angela.dai@tum.de

**Abstract.** Object recognition has seen significant progress in the image domain, with focus primarily on 2D perception. We propose to leverage existing large-scale datasets of 3D models to understand the underlying 3D structure of objects seen in an image by constructing a CAD-based representation of the objects and their poses. We present Mask2CAD, which jointly detects objects in real-world images and for each detected object, optimizes for the most similar CAD model and its pose. We construct a joint embedding space between the detected regions of an image corresponding to an object and 3D CAD models, enabling retrieval of CAD models for an input RGB image. This produces a clean, lightweight representation of the objects in an image; this CAD-based representation ensures a valid, efficient shape representation for applications such as content creation or interactive scenarios, and makes a step towards understanding the transformation of real-world imagery to a synthetic domain. Experiments on real-world images from Pix3D demonstrate the advantage of our approach in comparison to state of the art. To facilitate future research, we additionally propose a new image-to-3D baseline on ScanNet which features larger shape diversity, real-world occlusions, and challenging image views.

## 1 Introduction

Object recognition and localization in images has been a core task of computer vision with a well-studied history. Recent years have shown incredible progress in identifying objects in RGB images by predicting their bounding boxes or segmentation masks [9,16,26]. Although these advances are very promising, recognizing 3D attributes of objects such as shape and pose is crucial to many real-world applications. In fact, 3D perception is fundamental towards human understanding of imagery and real-world environments – from a single RGB image a human can easily perceive geometric structure, and is paramount for enabling higher-level scene understanding such as inter-object relationships, or interaction with an environment by exploration or manipulation of objects.

At the same time, we are now seeing a variety of advances in understanding the shape of a single object from image view(s), driven by exploration of various geometric representations: voxels [6,40,43], points [10,44], meshes [8,14,41], and implicit surfaces [30,32]. While these generative approaches have shown significant promise in

inferring the geometry of single objects, these approaches tend to generate geometry that may not necessarily represent a valid shape, with tendency towards noise or over-smoothing, and excessive tessellation. Such limitations render these results unsuitable for many applications, for instance content creation, real-time robotics scenarios, or interaction in mixed reality environments. In addition, the ability to digitize the objects of real world images to CAD models opens up new possibilities in helping to bridge the real-synthetic domain gap by transforming real-world images to a synthetic representation where far more training data is available.

In contrast, we propose Mask2CAD to join together the capabilities of 2D recognition and 3D reconstruction by leveraging CAD model representations of objects. Such CAD models are now readily available [3] and represent valid real-world object shapes, in a clean, compact representation – a representation widely used by existing production applications. Thus, we aim to infer from a single RGB image object detection in the image as well as 3D representations of each detected object as CAD models aligned to the image view. This provides a geometrically clean, compact reconstruction of the objects in an image, and a lightweight representation for downstream applications.
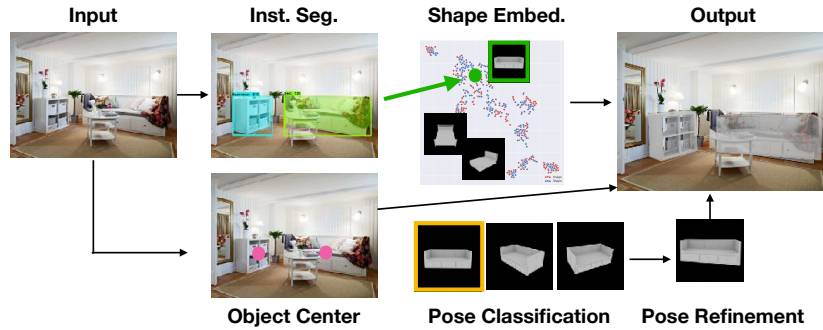


**Fig. 1.** Mask2CAD aims to predict object mask, and 3D shape in the scene. We achieve this by formulating an image-shape embedding learning problem. Combined with pose and object center prediction, Mask2CAD outputs realistic 3D shapes of objects from a single RGB image input. The entire system is differentiable and learned end-to-end.

Our Mask2CAD approach jointly detects object regions in an image and learns to map these image regions and CAD models to a shared embedding space (See Figure 1). At train time we learn a joint embedding which brings together corresponding image-CAD pairs, and pushes apart other pairs. At test time, we retrieve shapes by their renderings from the embedding space. To align the shapes to the image, we develop a pose prediction branch to classify and refine the shape alignment. We train our approach on the Pix3D dataset [38], achieving more accurate reconstructions than state of the art. Importantly, our retrieval-based approach allows adaptation to new domain by simply adding CAD models to the CAD model set without any re-training. Experiments on unseen shapes of the Pix3D dataset [38] show notable improvement when we have access to all CAD models at test time (but no access to corresponding RGB images of the

unseen models and no re-training). By leveraging CAD models as shape representation, we are able to predict multiple distinct 3D objects per image efficiently (approximately 60ms per image).

In addition to Pix3D, we also apply Mask2CAD on ScanNet and propose the first single-image to 3D object reconstruction baseline. Compared to Pix3D, this dataset contains 25K images, an order of magnitude more 3D shapes, complex real-world occlusions, diverse views and lighting conditions. Despite these challenges, Mask2CAD still manages to place appropriate CAD models that match the image observation (see Figure 5). We hope Mask2CAD could serve as a benchmark for future retrieval methods and reference for generative methods.

Mask2CAD opens up possibilities for object-based 3D understanding of images for content creation and interactive scenarios, and provides an initial step towards transforming real images to a synthetic representation.
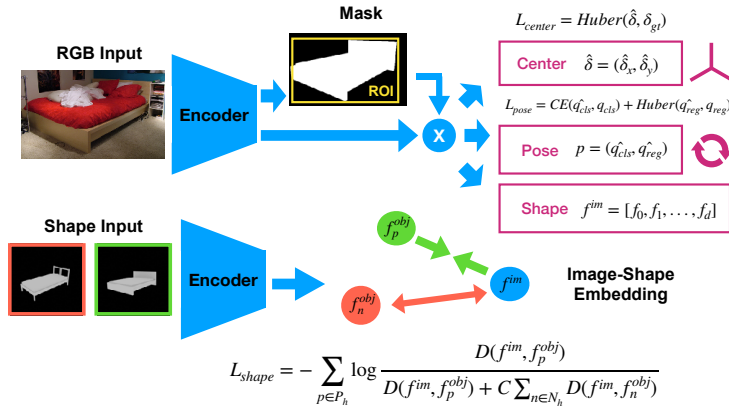


**Fig. 2.** Overview of our Mask2CAD approach for joint object segmentation and shape retrieval from a single RGB image. At train time, object detection is performed on an RGB image to produce a bounding box, segmentation mask, and feature descriptor for each detected object. The object feature descriptor is then used to train for an image-CAD embedding space for shape retrieval, as well as pose regression for the object rotation and center regression for its location. The embedding space is constructed through a triplet loss with corresponding and non-corresponding shapes to a detected object region of an image.

## 2 Related work

**Object Recognition in Images.** Our work draws inspiration from the success of 2D object detection and segmentation in the image domain, where myriad methods have been developed to predict 2D object bounding boxes and class labels from a single image input [11,26,28,33,34]. Recent approaches have focused on additionally predicting instance masks for each object [16,22]. We build from this 2D object detection and segmentation to jointly learn to predict shape as well.

**Single-View Object Reconstruction.** In recent years, a variety of approaches have been developed to infer 3D shape from a single RGB image observation, largely focusing on the single object scenario and exploring a variety of shape representations in the context of learning-based methods. Regular voxel grids are a natural extension of the 2D image domain, and have been shown to effectively predict global shape structures [6,43], but remain limited by computational and memory constraints when scaling to high resolutions, as well as uneconomical in densely representing free space. Thus methods which focus representation power on surface geometry have been developed, including hierarchical approaches on voxels such as octrees [35,40], or point-based representations [10,44]. More recently, methods have been developed to predict triangle meshes, largely based on strong topological assumptions such as deforming an existing template mesh [41] while free-form generative approaches tend to remain limited by computational complexity to very small numbers of vertices [8]. Implicit reconstruction approaches have also shown impressive shape reconstruction results at relatively high resolution by predicting the occupancy [30] or signed distance field value [32] for point sampled locations.

Mesh R-CNN [12] pioneered an approach to extend such single object reconstruction to jointly detect and reconstruct shape geometry for each detected object in an RGB image. Mesh R-CNN extends upon the object recognition pipeline of Mask R-CNN [16] to predict initial voxel-based occupancy of an object, which is then refined by a graph convolutional network to produce an output mesh for each object.

While these approaches for object reconstruction have shown significant promise in predicting general, structural shape properties, due to the low-level nature of the reconstruction approaches (generating on a per-voxel/per-point basis), the reconstructed objects tend to be noisy or oversmoothed, may not represent valid real-world shapes (e.g., disconnected in thin regions, missing object symmetries), and inefficiently represented in geometry (e.g., over-tessellated to achieve higher resolutions). In contrast, our Mask2CAD approach leverages existing CAD models of objects to jointly segment and retrieve the 3D shape for each object in an image, producing both an accurate reconstruction and clean, compact geometric representation.

**CAD Model Priors for 3D Reconstruction.** Leveraging geometric model-based priors for visual understanding has been established near the inception of computer vision and robotic understanding [2,5,36], although constrained by the geometric models available. With increasing availability of larger-scale CAD model datasets [3,38], we have seen a rejuvenation in understanding the objects of a scene by retrieving and aligning similar CAD models. Most methods focus on aligning CAD models to RGB-D scan, point cloud geometry, or 2D-3D surface mapping though various geometric feature matching techniques [1,13,19,21,24,37]. Izadinia and Seitz [20] and Huang et al. [17] propose to optimize for both scene layouts and CAD models of objects from image input, leveraging analysis-by-synthesis approaches; these methods involve costly optimization for each input image (minutes to hours).

From single image views of a object, Li et al. [25] propose a method to construct a joint embedding space between RGB images and CAD models by first constructing the space based on handcrafted shape similarity descriptors, and then optimizing for the image embeddings into the shape space. Our approach also optimizes for a joint

embedding space between image views and CAD models in order to perform retrieval; however, we construct our space by jointly learning from both image and CAD in an end-to-end fashion without any explicit encoding of shape similarity.

## 3   Mask2CAD

### 3.1   Overview

From a single RGB image, Mask2CAD detects and localizes objects by recognizing similar 3D models from a candidate set, and inferring their pose alignment to the image. We focus on real-world imagery and jointly learn the 2D-3D relations in an end-to-end fashion. This produces an object-based reconstruction and understanding of the image, where each object by nature is characterized by a valid, complete 3D model with a clean, efficient geometric representation.

Specifically, from an input image, we first detect all objects in the image domain by predicting their bounding boxes, class labels, and segmentation masks. From these detected image regions, we then learn to construct a shared embedding space between these image regions and 3D CAD models of objects, which enables retrieving a geometrically similar model for the image observation. We simultaneously predict the object alignment to the image as 5 dof pose optimization (z-depth translation given), yielding a 3D understanding of the objects in the image.

*Object Detection*  For object detection, we build upon ShapeMask [22], a state-of-the-art instance segmentation approach. ShapeMask takes as input an RGB image and outputs detected objects characterized by their bounding boxes, class labels, and segmentation masks. The one-stage detection approach of RetinaNet [26] is leveraged to generate object bounding box detections, which are then refined into instance masks by a learned set of shape priors. We modify it to leverage the learned features for our 3D shape prediction. Each bounding box detection is used to crop features from the corresponding level of feature pyramid to produce a feature descriptor $F_i$ for the instance mask prediction $M_i$ (e.g. 32x32) of object $i$; we then take the product $M_i \circ F_i$ as the representative feature map for object $i$ as seen in the image. This is then used to inform the following CAD model retrieval and pose alignment.

### 3.2   Joint Embedding Space for Image-CAD Retrieval

The core of our approach lies in learning to seamlessly map between image views of an object and 3D CAD models, giving an association between image and 3D geometry. The CAD models represent an explicit prior on object geometry, providing an inherently clean, complete, and compact 3D representation of an object. We learn this mapping between image-CAD by constructing a shared embedding space – importantly, as we show in Section 4, our approach to jointly learn this embedding space constructs a space that is robust to new, unseen CAD models at test time.

Constructing a joint embedding space between image regions and 3D object geometry requires mapping across two very different domains, where in contrast to a geometric CAD model, an image is view-dependent and composed of the interaction of scene

geometry with lighting and material. To facilitate the construction of this shared space between, we thus represent each object similar to a light field descriptor [4], rendering a set of $k$ views $I_0^i, ..., I_k^i$ for an object $O_i$. For all our experiments, we use $k = 16$; the set of canonical views for each object is determined by K-medoid clustering of the views seen of the object category in the training set. In addition, we augment the pool of anchor-positive pairs by using slightly jittered groundtruth views of the objects.

The embedding space is then established between the image region features $M_i \circ F_i$ from the detection, and the 3D object renderings $I_0^j, ..., I_k^j$. Representative features for the image region descriptions and object renderings are extracted by a series of convolutional layers applied to each input. The convolutional networks to extract these features are structured symmetrically, although we do not share weights due to the different input domains (See Sec. 3.4 for more details). We refer to the resulting extracted feature descriptors as $f^{im}$ and $f^{obj}$ for the image regions and object views, respectively. The $f^{im}$ come from the regions of interest (ROI) shared with the 2D detection and segmentation branch. More specifically, the encoder backbone is a ResNet feature pyramid network and the decoder is a stack of 3x3 convolution layers on the ROI features.

We guide the construction of the embedding space with a noise contrastive estimation loss [31] for $f^{im}$ describing a detected image region

$$L_c = -\sum_{p \in P_h} \log \frac{D(f^{im}, f_p^{obj})}{D(f^{im}, f_p^{obj}) + C \sum_{n \in N_h} D(f^{im}, f_n^{obj})} \tag{1}$$

where $f_p^{obj}$ represents the feature descriptor of a corresponding 3D object to the image region, $f_n^{obj}$ the feature descriptor of a non-corresponding object, $C$ a weighting parameter, and $D$ the cosine distance function:

$$D(x, y) := \frac{1}{\tau}(\frac{x}{||x||})^T(\frac{y}{||y||}) \tag{2}$$

where $\tau$ is the temperature. $P_h$ and $N_h$ denote the set of hard positive and negative examples for the image region. Details of hard-example mining are provided in the next section. The positively corresponding objects are determined by the CAD annotations to the images, and negatively corresponding objects are the non-corresponding CAD renderings in the training batch.

Since the object detection already provides class of the object, the negatives are only sampled from the shapes under the same class; that is, our embedding spaces are constructed for each class category although the weights are shared among them.

Empirically, we find it important to place more sampling weights on the rare classes because the number of valid pairs scale *quadratically* with the number of same-class examples in the batch. We apply the inverse square root method as in [15] to enhance the rare class examples with a threshold $t = 0.1$, which leads to improved performance on rare classes without compromising dominant classes.

**Hard example mining.** Hard example mining is known to be crucial for embedding learning, as most examples are easy and do not contain much information to improve the model. We employ both hard positive and hard negative mining in Mask2CAD as

follows. For each image region (anchor), we sample top-$P_h$ positive object views and top-$N_h$ negative object views by their distances to anchor. Similar to [22], we sample $Q$ objects for each image during the training ($Q = 8$). Since the number of objects in a batch scales linearly with $Q$, we set $P_h = 4Q = 32$ and $N_h = 16Q = 128$. Summation over hard examples allows the loss to focus on difficult cases and perform better.

**Shape Retrieval.** Once this embedding space is constructed, we can then perform shape retrieval at test time to provide a 3D understanding of the objects in an image. An input image at test time is processed by the 2D detection to provide a bounding box, segmentation mask, and feature descriptor for each detected object. We then use a nearest neighbor retrieval into the embedding space with $N_k = 1$ based on cosine distance to find the most similar CAD model for each detected object. We have tried larger $N_k$ values and majority vote schemes but did not see any performance gain.

### 3.3  Pose Prediction

We additionally aim to predict the pose of the retrieved 3D object such that it aligns best to the input image. We thus propose a pose prediction branch which outputs the rotation and translation of the object. Starting with the $M_i \circ F_i$ feature map for a detected object, the object translation is directly regressed with a Huber loss [18] as follows:

$$L_\delta(x) = \begin{cases} \frac{1}{2}x^2 & \text{for } |x| \leq \delta, \\ \delta(|x| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases} \tag{3}$$

The object rotation is simultaneously predicted; the rotation is first classified to a set of $K$ discretized rotation bins using cross entropy loss; this coarse estimate is then refined through a regression step using a Huber loss. This coarse-to-fine approach helps to navigate the non-euclidean rotation space, and enables continuous rotation predictions.

For rotation prediction, we represent the rotation as a quaternion, and compute the set of rotation bins by K-medoid clustering based on train object rotations. To further refine this coarse prediction, we then predict a refined rotation by estimating the delta from the classified bin using a Huber loss. The delta is represented as a $R^4$ quaternion. We initialize the bias of the last layer with $(0.95, 0, 0, 0)$ such that the quaternion is close to identity transform at the beginning. Note that during training, we only train the refinement for classified rotations within $\theta$ to avoid regressing to dissimilar targets.

To obtain full prediction in the camera space, we need to predict the translation of the object in addition to the shape and rotation. A naive approach is to use the bounding box center as the object center in 2D and cast a ray through the center to intersects with the given groundtruth z-plane. Unfortunately, the bounding boxes tend to be unstable against the rotation and their centers can end up far from the actual object center.

We thus regress the object center as a bounding regression problem. More specifically, for each ROI, we task the network with predicting $(\delta_x, \delta_y)$, where the $\delta$s are the shift between bounding box center and actual object center as a ratio of object width and height. At train time, we optimize $(\delta_x, \delta_y)$ with the aforementioned Huber loss. At test time, we apply $(\delta_x, \delta_y)$ to the box center to obtain the object 3D translations (assuming groundtruth depth is given [12]).

### 3.4   Implementation details

We use ShapeMask [22] as the instance segmentation backbone. The model backbone is initialized from COCO-pretrained checkpoint and uses ResNet-50 architecture so as to be comparable to Mesh R-CNN in our experiments. The shape rendering branch uses a lightweight ResNet-18 backbone initialized randomly.

   We freeze the weights of the backbone ResNet-50 layers after initialization and optimize both branches jointly for 48K iterations until convergence (about 1000 epochs for Pix3D), which takes approximately 13 hours. The learning rate is set to be 0.08 and decays by 0.1 at 32K and 40K iterations. The losses for the retrieval and pose estimation are weighted with 0.5, 0.25, and 5.0 for the embedding loss, pose classification loss, and pose regression loss. We use $C = 1.5$ and $\tau = 0.15$ in our contrastive loss, and Huber loss margin of $\delta = 0.15$ for the pose and center regression. For pose prediction, we set $K = 16$ bins, and $\theta = \pi/6$.

   For each example, we randomly sample 3 out of $k = 16$ canonical view renderings and one jittered groundtruth view rendering to add to the contrastive learning pool. Similar to ShapeMask, we apply ROI jittering to the image region for training the segmentation, embedding, and pose estimation branches. The noise is set to 0.025 following ShapeMask. We also apply data augmentation by horizontal image flips with 50% probability. For such image flips, the pose labels were also adjusted accordingly.

## 4   Experiments

We evaluate our approach on the Pix3D dataset [38], which comprises $10,069$ images annotated with corresponding 3D models of the objects in the images. We aim to jointly detect and predict the 3D shapes for the objects in the images. We evaluate on the train/test split used by Mesh R-CNN [12] for the same task. Additionally, we propose the first single-image 3D object reconstruction baseline on the ScanNet dataset [7], which tends to contain more cluttered, in-the-wild views of objects.

*Evaluation metric.*   We adopt the popular metrics from 2D object recognition, and similar to Mesh R-CNN [12], evaluate $AP^{box}$ and $AP^{mask}$ on the 2D detections, and $AP^{mesh}$ on the 3D shape predictions for the objects. Similar to Mesh R-CNN, we evaluate $AP^{mesh}$ using the precision-recall for $F1^{0.3}$. However, note that while Mesh R-CNN evaluate these metrics at IoU 0.5 (AP50), we adopt the standard COCO object detection protocol of AP50-AP95 (denoted as AP), averaging over 10 IoU thresholds of $0.5 : 0.05 : 0.95$ [27]. This enables characterization of high-accuracy shape reconstructions captured at more strict IoU thresholds, demonstrating a more comprehensive description of the accuracy of the shape predictions. In addition to AP, we also report individual $AP^{mesh}$ scores for IoU thresholds of 0.5 and 0.75 following Mask R-CNN [16]. For better reproducibility, we report every metric as an average of 5 independent runs throughout this paper.

**Comparison to state of the art.** We compare our Mask2CAD approach for 3D object understanding from RGB images by joint segmentation and retrieval to Mesh R-CNN [12], who first propose this task on Pix3D [38]. Table 1 shows our shape prediction

**Table 1.** Performance on Pix3D [38] $\mathcal{S}_1$. We report mean AP$^{\text{mesh}}$ as well as per category AP$^{\text{mesh}}$. AP is averaged from AP50-AP95 following the COCO detection protocol. All AP performances are in %. We outperform the state-of-the-art approach on all AP metrics. This improvement mostly derives from maintaining more robust performance in the high AP regime above AP50. Additionally, we observe that Mask2CAD performs well on furniture categories and not so well on tools and miscellaneous objects which exhibit highly irregular shapes.

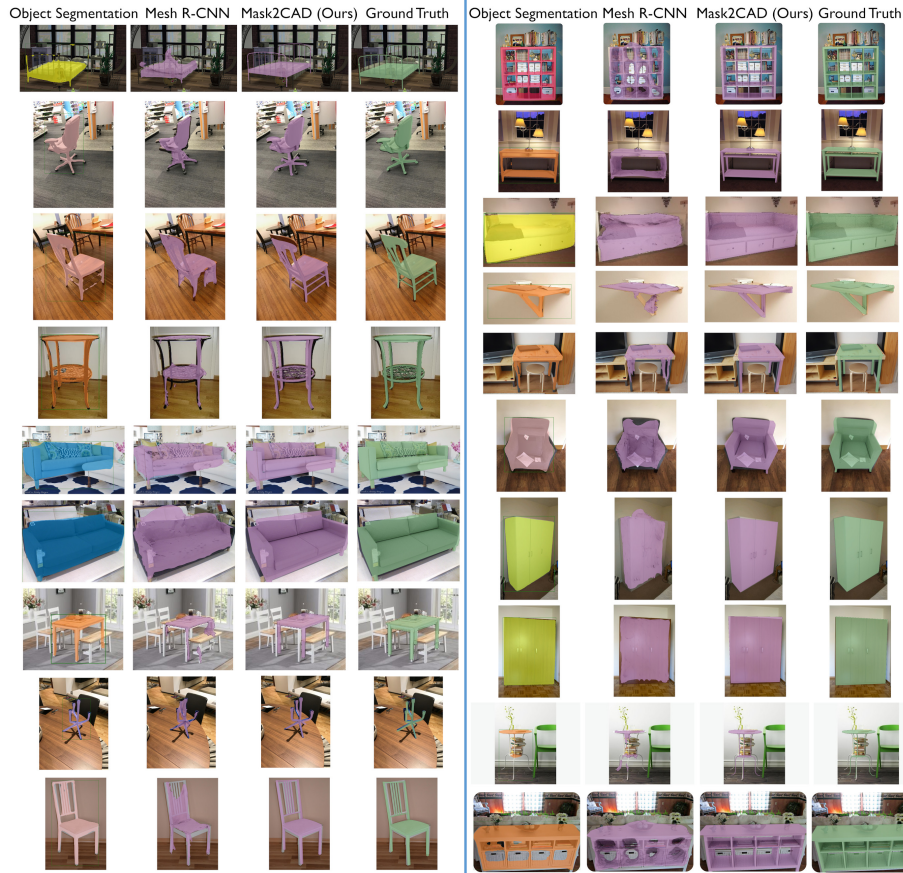| Pix3D $\mathcal{S}_1$ | AP | AP50 | AP75 | chair | sofa | table | bed | desk | bkcs | wrdrb | tool | misc |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mesh R-CNN [12] | 17.2 | 51.2 | 7.4 | 17.6 | 30.0 | 11.0 | 20.0 | 21.0 | 10.1 | 14.3 | **6.5** | **24.5** |
| Mask2CAD | **33.2** | **54.9** | **30.8** | **19.6** | **55.8** | **29.2** | **39.4** | **31.6** | **42.4** | **60.3** | 4.2 | 15.9 |



**Fig. 3.** Mask2CAD predictions on Pix3D [38]. The detected object is highlighted on the lefthand side of each column, with shape predictions denoted in purple, and ground truth in green. In contrast to Mesh R-CNN [12], our approach can achieve more accurate shape predictions with geometry in a clean, efficient representation.

results in comparison to Mesh-RCNN on their $\mathcal{S}_1$ split of the Pix3D dataset. We evaluate $AP^{mesh}$, averaged over all class categories, as well as per-category. In contrast to Mesh R-CNN, whose results show effective coarse predictions but suffer significantly at AP75, our shape and pose predictions maintain high-accuracy reconstructions. We show qualitative comparisons in Figure 3.

**Table 2.** Performance on Pix3D [38] with ground truth object bounding boxes given. We report Chamfer distance, Normal consistency and F1 scores. Note that for these experiments, the Mesh R-CNN-based approaches are additionally provided the ground truth scale in the depth dimension of the object.

| Pix3D $\mathcal{S}_1$ gt | Chamfer ↓ | Normal ↑ | $F1^{0.1}$ ↑ | $F1^{0.3}$ ↑ | $F1^{0.5}$ ↑ |
|---|---|---|---|---|---|
| Mask R-CNN + Pixel2Mesh [12] | 1.30 | 0.70 | 16.4 | 51.0 | 68.4 |
| Mesh R-CNN (Voxel-Only) [12] | 1.28 | 0.57 | 9.9 | 37.3 | 56.1 |
| Mesh R-CNN (Sphere-Init) [12] | 1.30 | 0.69 | 16.8 | 51.4 | 68.8 |
| Mesh R-CNN [12] | 1.11 | 0.71 | 18.7 | 56.4 | 73.5 |
| Mask2CAD (Ours) | **0.99** | **0.74** | **25.6** | **66.4** | **79.3** |

Additionally, we compare to several state-of-the-art single object reconstruction approaches on Pix3D $\mathcal{S}_1$ in Table 2; for each approach we provide ground truth 2D object detections, i.e. perfect bounding boxes. We evaluate various characteristics of the shape reconstruction. We also evaluate the Chamfer distance, normal consistency, and F1 at thresholds 0.1, 0.3, 0.5, using randomly sampled points on the predicted and ground truth meshes, where meshes are scaled such that the longest edge of the ground-truth meshs bounding box has length 10. Chamfer distance and normal consistency provide more global measures of shape consistency with the ground truth, but can tend towards favoring averaging, while F1 scores tend to be more robust towards outliers, and F1 at lower thresholds in particular indicates the ability to predict highly-accurate shapes. Note the competing approaches have been provided the ground truth scale in the depth dimension at test time, while our approach directly retrieves it from the training set. Nonetheless, our approach can provide higher-accuracy predictions as seen in the F1 scores at 0.1 and 0.3.

**Implicit learning of shape similarity.** In Figure 4, we visualize our learned embedding space by t-SNE [29], for image regions and CAD models of the *sofa* and *bookcase* class categories (we refer to the supplemental material for additional visualizations of the learned embedding spaces). We find that not only do the images and shapes mix together in this embedding space, despite that it is constructed without any knowledge of shape similarity – only image-CAD associations –, geometrically similar shapes tend to cluster together.

**Can the image-shape embedding space generalize to new 3D models?** Our joint image-CAD model embedding space constructed during train time leverages ground truth annotations of CAD models to images, which can be costly to acquire. During

**Fig. 4.** t-SNE embeddings of Mask2CAD for the sofa (top) and bookcase (bottom) classes. More visualizations can be found in the Supp. Materials. Red points correspond to images, and blue to shapes. Both images and shapes mix well together in the embedding space. Note that despite lack of shape similarity information during training, clusters tend to form together in geometric similarity, e.g., L-shaped sofas (yellow), single seat sofas without armrests (orange), single seat sofas with armrests (blue), double seat sofas (green). This stands in contrast to the embedding space construction of [25] which explicitly enforces shape similarity in its light field descriptors.

inference time, however, we can still embed new 3D models into the space without training, by using our trained model to compute their feature embeddings. Our embedding approach generalizes well in incorporating these new models.

We demonstrate this on the $\mathcal{S}_2$ split of Pix3D, training on a subset of the 3D CAD models, with test images comprising views of a disjoint set of objects than those in the training set. Generalization under this regime is difficult, particularly for a retrieval-based approach. However, in Table 3 we show clear improvements when using all available CAD models at test time in comparison to only the CAD models in the train set, despite not having seen any of the new objects nor their corresponding image views.

To help the model generalize better, we apply more data augmentation than the $\mathcal{S}_1$ split, including HSV-space jittering, random crop and resize of the renderings, and augmenting the box and image jittering magnitude as used in ShapeMask [22].

**Table 3.** Test-time generalization on Pix3D [38] $\mathcal{S}_2$. The performance improves on all categories with the addition 139 of CAD models at test time without re-training.

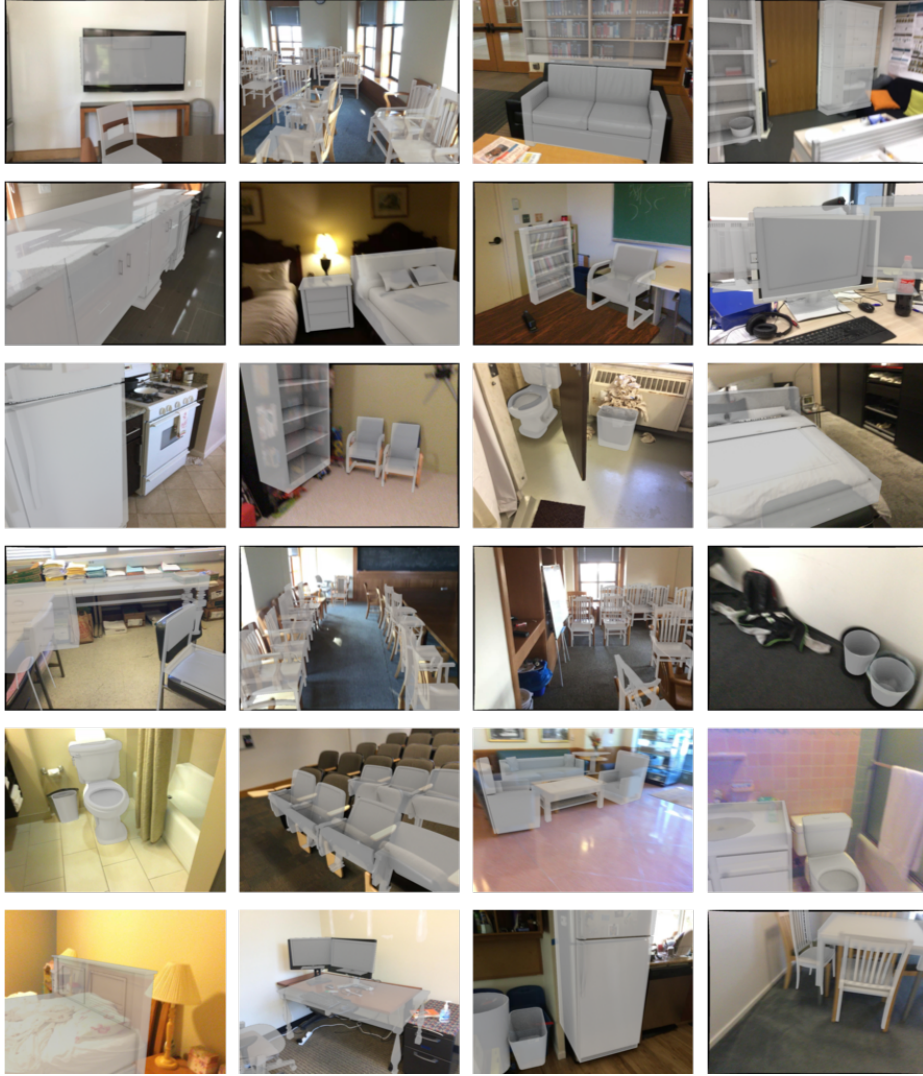| Pix3D $\mathcal{S}_2$ | AP | AP50 | AP75 | chair | sofa | table | bed | desk | bkcs | wrdrb | tool | misc |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mask2CAD (Ours) | 6.5 | 17.3 | 3.8 | 3.2 | 35.4 | 1.2 | 14.0 | 0.2 | 2.2 | 1.6 | 0.6 | 0.0 |
| Mask2CAD (Ours) + CAD | **8.2** | **20.7** | **4.8** | **4.5** | **37.8** | **3.6** | **16.9** | **2.7** | **2.2** | **5.3** | **0.9** | **0.1** |



**Fig. 5.** Example Mask2CAD predictions on ScanNet [7] images. Our approach shows encouraging results in its application to the more diverse set of image views, lighting, occlusions, and object categories of ScanNet.

**Comparison with ShapeNet reconstruction methods.** In Table 4, we compare the Pix3D $\mathcal{S}_1$ model on the validation set with the other methods that train on ShapeNet with real data augmentation. The evaluation protocol and implementation follows [38]. Mask2CAD results are reported on 1165 chairs in the $\mathcal{S}_1$ test split of Mesh R-CNN [12], as an average over 5 independent runs. Surprisingly, Mask2CAD achieves significantly better shape predictions than the state-of-the-art methods (0.288 IoU, and 0.013 Chamfer Distance), showing the capability of retrieval.

**Table 4.** Mask2CAD on Pix3D [38] $\mathcal{S}_1$ test split in comparison with other methods that train on ShapeNet models with real data augmentation.

|         | Mask2CAD (Ours) | FroDO [23] | Sun et al. [38] | MarrNet [42] | 3D-R2N2 [6] |
|---------|:---------------:|:----------:|:---------------:|:------------:|:-----------:|
| IoU     | **0.613**       | 0.325      | 0.282           | 0.231        | 0.136       |
| Chamfer | **0.086**       | 0.099      | 0.119           | 0.144        | 0.239       |

**Baseline on ScanNet Dataset.** We additionally apply Mask2CAD to real-world images from the ScanNet dataset [7], which contains RGB-D video data of 1513 indoor scenes. We use the 25K frame subset provided by the dataset for training and validation. The train/val split contains 19387/5436 images respectively, and the images come from separate scenes with distinct objects. Compared to Pix3D, this dataset has an order of magnitude more shapes, as well as many more occlusions, diverse image views, lighting conditions, and importantly, metric 3D groundtruth of the scene. We believe this could be a suitable benchmark for object 3D prediction from a real single image.

We use the CAD labels from Scan2CAD [1] by projecting the CAD models to each image view and use the amodal box, mask, pose, and shape for training. We additionally remove the objects whose centers are out of frame from training and evaluation. We also remove the object categories that appear less than 1000 times in the training split, resulting in eight categories: bed, sofa, chair (inc. toilet), bin, cabinet (inc. fridge), display, table, and bookcase. Regarding shape similarity, we adopt F score = 0.5 as the threshold for Mesh AP computation, because the Scan2CAD annotations come from ShapeNet and do not provide exact matches to the images. As Scan2CAD provides 9-DoF annotation for each object, we apply the groundtruth z depth and (x, y, z) scale to the predicted shape before computing the shape similarity metrics. We trained Mask2CAD for 72000 iterations with HSV-color, ROI, and image scale jittering using the same learning rate schedule as Pix3D. The quantitative results are reported in Table 5. Despite the complexity of ScanNet data, Mask2CAD manages to recognize the object shapes in these images, as shown in Figure 5. Our CAD model retrieval and alignment shows promising results and a potential for facilitating content creation pipelines.

**Runtime.** At test time Mask2CAD is efficient and runs at approximately 60 milliseconds per 640 by 640 image on Pix3D, including 2D detection and segmentation as well as shape retrieval and pose estimation.

**Table 5.** Performance on ScanNet [7]. We report mean $AP^{mesh}$ as well as per category $AP^{mesh}$ following Pix3D protocol.

| ScanNet 25K | AP | AP50 | AP75 | bed | sofa | chair | cab | bin | disp | table | bkcs |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Mask2CAD(Ours) | 8.4 | 23.1 | 4.9 | 14.2 | 13.0 | 13.2 | 7.5 | 7.8 | 5.9 | 2.9 | 3.1 |

*Limitations.* While Mask2CAD shows promising progress in attaining 3D understanding of the objects from a single image, we believe there are many avenues for further development. For instance, our retrieval-based approach can suffer in the case of objects that differ too strongly from the existing CAD model database, and we believe that mesh-based approaches to deform and refine geometry [39,41] have significant potential to complement our approach. Additionally, we believe a holistic 3D scene understanding characterizing not only the objects in an environment but all elements in the scene is a promising direction for 3D perception and semantic understanding.

## 5   Conclusion

We propose Mask2CAD, a novel approach for 3D perception from 2D images. Our method leverages a CAD model representation, and jointly detects objects for an input image and retrieves and aligns a similar CAD model to the detected region. We show that our approach produces accurate shape reconstructions and is capable of generalizing to unseen 3D objects at test time. The final output of Mask2CAD is a CAD-based object understanding of the image, where each object is represented in a clean, lightweight fashion. We believe that this makes a promising step in 3D perception from images as well as transforming real-world imagery to a synthetic representation, opening up new possibilities for digitization of real-world environments for applications such as content creation or domain transfer.

## Acknowledgments

# References

1. Avetisyan, A., Dahnert, M., Dai, A., Savva, M., Chang, A.X., Nießner, M.: Scan2cad: Learning cad model alignment in rgb-d scans. CVPR (2019)
2. Binford, T.O.: Survey of model-based image analysis systems. The International Journal of Robotics Research **1**(1), 18–64 (1982)
3. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., Xiao, J., Yi, L., Yu, F.: ShapeNet: An Information-Rich 3D Model Repository. Tech. Rep. arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago (2015)
4. Chen, D.Y., Tian, X.P., Shen, Y.T., Ouhyoung, M.: On visual similarity based 3d model retrieval. In: Computer graphics forum. vol. 22, pp. 223–232. Wiley Online Library (2003)
5. Chin, R.T., Dyer, C.R.: Model-based recognition in robot vision. ACM Computing Surveys (CSUR) **18**(1), 67–108 (1986)
6. Choy, C.B., Xu, D., Gwak, J., Chen, K., Savarese, S.: 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In: European conference on computer vision. pp. 628–644. Springer (2016)
7. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: ScanNet: Richly-annotated 3D reconstructions of indoor scenes. In: Proc. Computer Vision and Pattern Recognition (CVPR), IEEE (2017)
8. Dai, A., Nießner, M.: Scan2mesh: From unstructured range scans to 3d meshes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5574–5583 (2019)
9. Duan, K., Bai, S., Xie, L., Qi, H., Huang, Q., Tian, Q.: Centernet: Keypoint triplets for object detection. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 6569–6578 (2019)
10. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 605–613 (2017)
11. Girshick, R.: Fast r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 1440–1448 (2015)
12. Gkioxari, G., Malik, J., Johnson, J.: Mesh r-cnn. arXiv preprint arXiv:1906.02739 (2019)
13. Grabner, A., Roth, P.M., Lepetit, V.: Location field descriptors: Single image 3d model retrieval in the wild. In: 2019 International Conference on 3D Vision (3DV). pp. 583–593. IEEE (2019)
14. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: A papier-mâché approach to learning 3d surface generation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 216–224 (2018)
15. Gupta, A., Dollar, P., Girshick, R.: Lvis: A dataset for large vocabulary instance segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5356–5364 (2019)
16. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Computer Vision (ICCV), 2017 IEEE International Conference on. pp. 2980–2988. IEEE (2017)
17. Huang, S., Qi, S., Zhu, Y., Xiao, Y., Xu, Y., Zhu, S.C.: Holistic 3D scene parsing and reconstruction from a single RGB image. In: European Conference on Computer Vision. pp. 194–211. Springer (2018)
18. Huber, P.J.: Robust estimation of a location parameter. In: Breakthroughs in statistics, pp. 492–518. Springer (1992)
19. Izadinia, H., Seitz, S.M.: Scene recomposition by learning-based icp. In: CVPR (2020)

20. Izadinia, H., Shan, Q., Seitz, S.M.: Im2cad. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5134–5143 (2017)
21. Kim, Y.M., Mitra, N.J., Huang, Q., Guibas, L.: Guided real-time scanning of indoor objects. In: Computer Graphics Forum. vol. 32, pp. 177–186. Wiley Online Library (2013)
22. Kuo, W., Angelova, A., Malik, J., Lin, T.Y.: Shapemask: Learning to segment novel objects by refining shape priors. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 9207–9216 (2019)
23. Li, K., Rünz, M., Tang, M., Ma, L., Kong, C., Schmidt, T., Reid, I., Agapito, L., Straub, J., Lovegrove, S., et al.: Frodo: From detections to 3d objects. arXiv preprint arXiv:2005.05125 (2020)
24. Li, Y., Dai, A., Guibas, L., Nießner, M.: Database-assisted object retrieval for real-time 3d reconstruction. In: Computer Graphics Forum. vol. 34, pp. 435–446. Wiley Online Library (2015)
25. Li, Y., Su, H., Qi, C.R., Fish, N., Cohen-Or, D., Guibas, L.J.: Joint embeddings of shapes and images via cnn image purification. ACM transactions on graphics (TOG) **34**(6), 1–12 (2015)
26. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: Proceedings of the IEEE international conference on computer vision. pp. 2980–2988 (2017)
27. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
28. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: Ssd: Single shot multibox detector. In: European conference on computer vision. pp. 21–37. Springer (2016)
29. Maaten, L.v.d., Hinton, G.: Visualizing data using t-sne. Journal of machine learning research **9**(Nov), 2579–2605 (2008)
30. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4460–4470 (2019)
31. Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748 (2018)
32. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 165–174 (2019)
33. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 779–788 (2016)
34. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: Advances in neural information processing systems. pp. 91–99 (2015)
35. Riegler, G., Osman Ulusoy, A., Geiger, A.: Octnet: Learning deep 3d representations at high resolutions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3577–3586 (2017)
36. Roberts, L.G.: Machine perception of three-dimensional solids. Ph.D. thesis, Massachusetts Institute of Technology (1963)
37. Shao, T., Xu, W., Zhou, K., Wang, J., Li, D., Guo, B.: An interactive approach to semantic modeling of indoor scenes with an rgbd camera. ACM Transactions on Graphics (TOG) **31**(6), 1–11 (2012)
38. Sun, X., Wu, J., Zhang, X., Zhang, Z., Zhang, C., Xue, T., Tenenbaum, J.B., Freeman, W.T.: Pix3D: Dataset and methods for single-image 3D shape modeling. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2974–2983 (2018)

39. Tang, J., Han, X., Pan, J., Jia, K., Tong, X.: A skeleton-bridged deep learning approach for generating meshes of complex topologies from single rgb images. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4541–4550 (2019)

40. Tatarchenko, M., Dosovitskiy, A., Brox, T.: Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2088–2096 (2017)

41. Wang, N., Zhang, Y., Li, Z., Fu, Y., Liu, W., Jiang, Y.G.: Pixel2mesh: Generating 3d mesh models from single rgb images. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 52–67 (2018)

42. Wu, J., Wang, Y., Xue, T., Sun, X., Freeman, B., Tenenbaum, J.: Marrnet: 3d shape reconstruction via 2.5 d sketches. In: Advances in neural information processing systems. pp. 540–550 (2017)

43. Wu, J., Zhang, C., Xue, T., Freeman, B., Tenenbaum, J.: Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In: Advances in neural information processing systems. pp. 82–90 (2016)

44. Yang, G., Huang, X., Hao, Z., Liu, M.Y., Belongie, S., Hariharan, B.: Pointflow: 3d point cloud generation with continuous normalizing flows. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 4541–4550 (2019)

## Appendix A: Additional Results on Pix3D

In Figure 7, we show more qualitative results of our approach on Pix3D [38]. Furthermore, we conduct ablation studies to shed light on the roles of each component in the system. Our analysis shows that shape, pose, and translation are all important for estimating the viewer-centered geometry, with shape retrieval having the most room for improvement, and box detection having the least. The analysis was done by replacing each predicted component with its ground truth counterpart. In terms of Mesh AP, groundtruth shapes help by +14.6, rotation by +10.2, and translation by +7.5. Surprisingly, groundtruth 2D boxes offer no improvement because the detections on Pix3D are very good ( 90 Box AP, similar to Mesh R-CNN) and the small advantage is offset by the distribution shift between train (jittered boxes) and test (perfect boxes) time. This agrees with what Mesh R-CNN reports, i.e. they also observed a loss when using ground truth boxes (6 point loss on Mesh AP50).

## Appendix B: Network Architecture Details

The Mask2CAD image-stream network architecture comprises 2D detection as bounding box, class label, and instance mask prediction, as well as our 3D shape retrieval and pose estimation. For the 2D detection, our architecture borrows from that of Shape-Mask [22]. For the 3D inference with shape embedding, pose classification, and pose regression, and object center prediction, these branches all use the same architecture as the coarse mask prediction branch of [22] (with the exception of the output layers). The inputs of these branches are the features from the region of interest (ROI) of detection backbone feature pyramid network. We detail each branch in Table 6, 7, and 8.

| Index | Inputs | Operation | Output shape |
|---|---|---|---|
| (1) | Input | Region of Interest (ROI) features | $32 \times 32 \times 256$ |
| (2) | (1) | $3\times$ (Conv($256 \rightarrow 256, 3 \times 3$) + BatchNorm + ReLU) | $32 \times 32 \times 256$ |
| (3) | (2) | (Conv($256 \rightarrow 256, 3 \times 3$) + BatchNorm + ReLU) | $32 \times 32 \times 128$ |
| (4) | (3) | AveragePool(axes=(0, 1)) | 128 |

**Table 6.** Network architecture of the shape embedding branch. The last convolution layer downsamples the number of channels from 256 to 128.

| Index | Inputs | Operation | Output shape |
|---|---|---|---|
| (1) | Input | Region of Interest (ROI) features | $32 \times 32 \times 256$ |
| (2) | (1) | $4\times$ (Conv($256 \rightarrow 256, 3 \times 3$) + BatchNorm + ReLU) | $32 \times 32 \times 256$ |
| (3) | (2) | AveragePool(axes=(0, 1)) | 256 |
| (4) | (3) | Linear($256 \rightarrow N_{pose} \times N_{class}$) | 160 |

**Table 7.** Network architecture of the pose prediction branch. For pose classification, the output is $N_{pose} = 16$ for each class $N_{class} = 10$. For the following pose regression after this classification, the architecture is identical except for using $N_{pose} = 4$ for predicting the regression quaternion instead of the 16 medoid bins.

| Index | Inputs | Operation | Output shape |
|---|---|---|---|
| (1) | Input | Region of Interest (ROI) features | $32 \times 32 \times 256$ |
| (2) | (1) | $4\times$ (Conv($256 \rightarrow 256$, $3 \times 3$) + BatchNorm + ReLU) | $32 \times 32 \times 256$ |
| (3) | (2) | AveragePool(axes=(0, 1)) | 256 |
| (4) | (3) | Linear($256 \rightarrow N_{center} \times N_{class}$) | 20 |

**Table 8.** Network architecture of the object center regression branch. The output is $N_{center} = 2$ for each class $N_{class} = 10$, where $N_{center}$ equals 2 for $(\delta_x, \delta_y)$.

## Appendix C: t-SNE visualizations for image-CAD embeddings

Figures 6, 8, 9 show the t-SNE visualizations of the image-shape embedding spaces for the bed, wardrobe, desk, table, tool, misc, and chair classes.
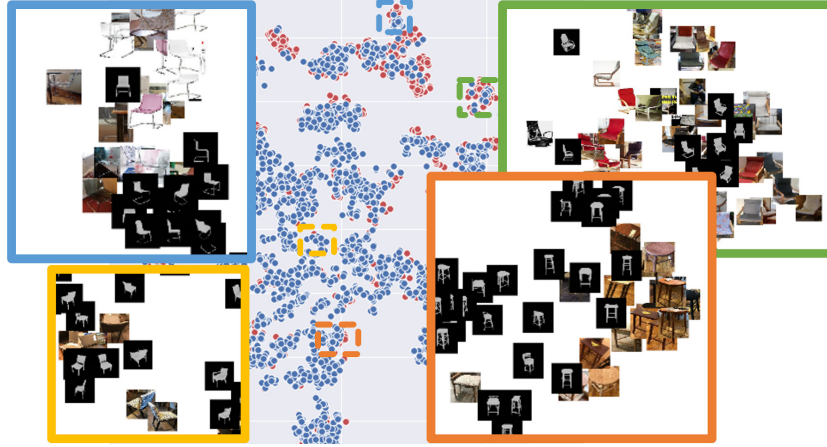


**Fig. 6.** t-SNE embedding of Mask2CAD for the chair class. Red points correspond to images, and blue to shapes.

**Fig. 7.** Additional qualitative results of Mask2CAD on Pix3D [38].
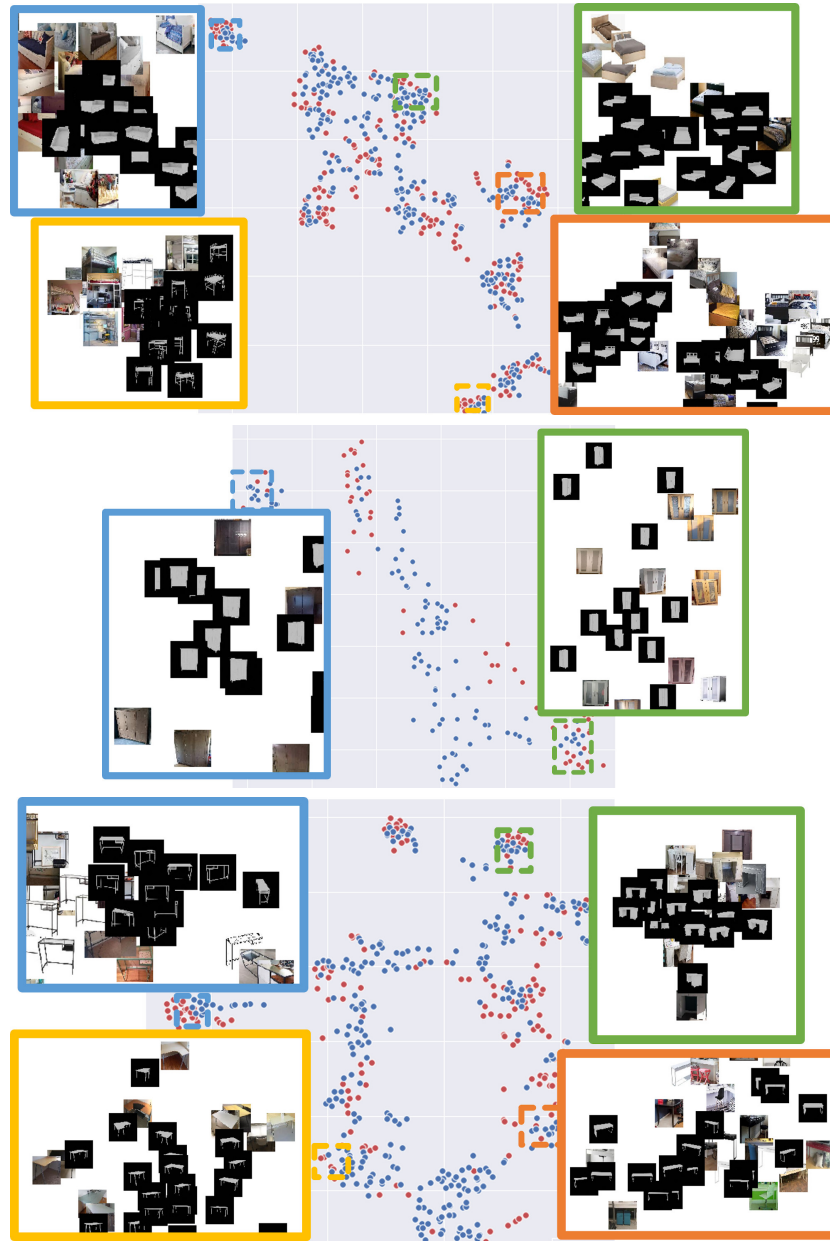
**Fig. 8.** t-SNE embeddings of Mask2CAD for the bed (top), wardrobe (middle) and desk (bottom) classes. Red points correspond to images, and blue to shapes.
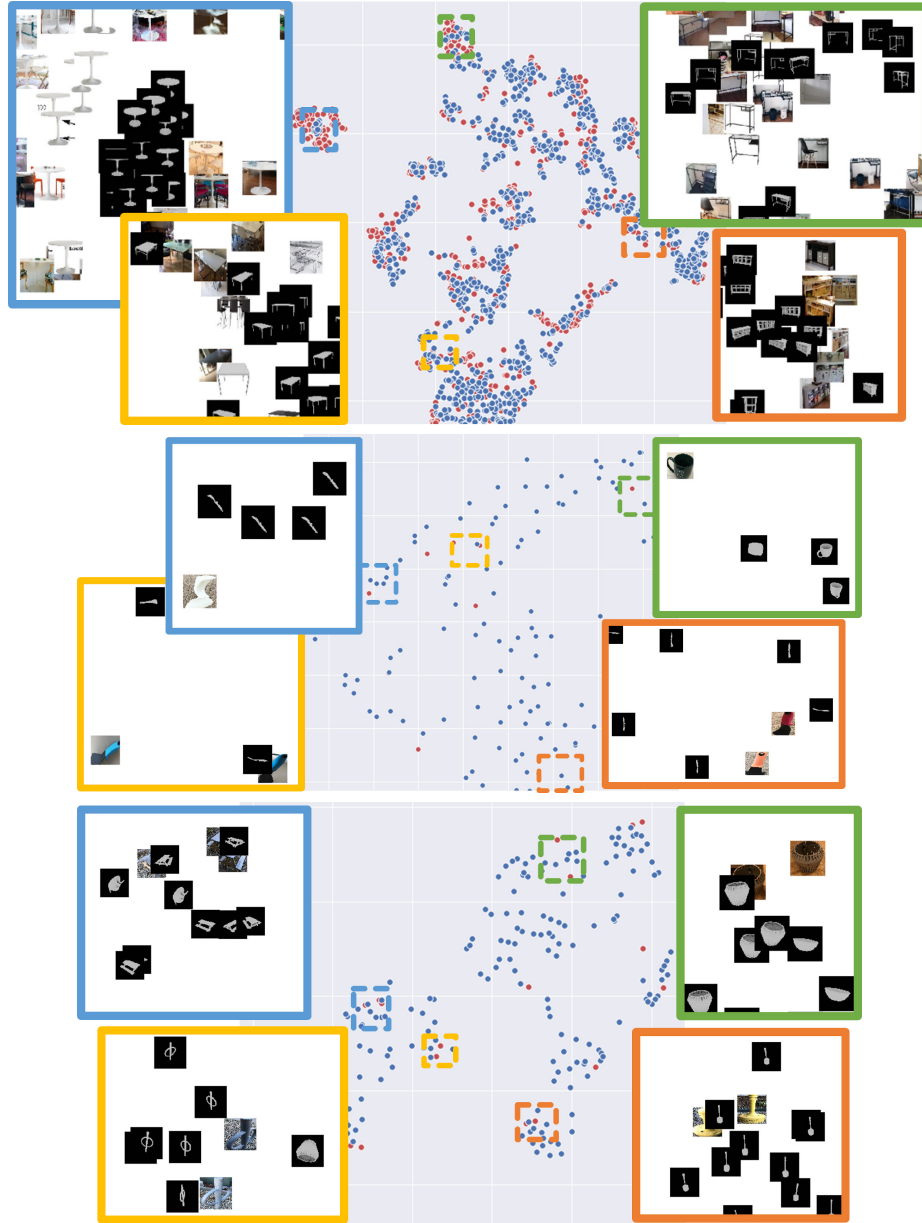
**Fig. 9.** t-SNE embeddings of Mask2CAD for the table (top), tool (middle) and misc (bottom) classes. Red points correspond to images, and blue to shapes.