

# CV Project Report (Spring2020)

## A Point Set Generation Network for 3D Object Reconstruction from a Single Image (CVPR December 2016)

### Paper by:

Haoqiang Fan\*  
Hao Su\*  
Leonidas Guibas

### Implemented by:

Team: Vicissitude  
Apoorva Srivastava (2019702014)  
Anurag Sahu (2018121004)  
Kajal Sanklecha (2019801006)  
Guide: Dr. Avinash Sharma  
TA: Aditya Aggarwal

# Contents

| S.No. | Topic                  | Page No. |
|-------|------------------------|----------|
| 1     | Contents               | 2        |
| 2     | Introduction           | 3        |
| 3     | Challenges             | 4        |
| 4     | Network Architecture   | 5        |
| 5     | Loss Function          | 8        |
| 6     | Implementation Details | 10       |
| 7     | Code and Results       | 11       |
| 8     | Future Tasks           | 15       |
| 9     | References             | 16       |

# Introduction

The current paper deals with the problem of 3D reconstruction from a single 2D image. Multiple techniques have been applied to achieve the same using Deep Learning and the output 3D are always either a collection of images or depth maps, or volumetric representations etc. However, these representations obscure the natural invariance of 3D shapes under geometric transformations and also suffer from a number of other issues like such representations lead to difficult trade offs between sampling resolution and net efficiency. Furthermore, they enshrine quantization artifacts that obscure natural invariances of the data under rigid motions etc.

In this paper, for the first time it has been tried to obtain the output as a collection of unordered 3D point coordinates.

There are following advantages of obtaining output in this form-

1. A point cloud is a simple, uniform structure that is easier to learn, as it does not have to encode multiple primitives or combinatorial connectivity patterns.
2. A point cloud allows simple manipulation when it comes to geometric transformations and deformations, as connectivity does not exist.
3. One advantage of point sets comes from its un-orderedness. Unlike 2D based representations like the depth map no topological constraint is put on the represented object.

But if point clouds are this much advantageous then the question arises why do the majority of 3D representations are either volumetric or collection of images? It is because extant deep net architectures for both discriminative and generative learning in the signal domain are well suited to data that is regularly sampled, such as images, audio, or video. However, most common 3D geometry representations, such as 2D meshes or point clouds are not regular structures and do not easily fit into architectures that exploit such regularity for weight sharing, etc. That is why the majority of extant works on using deep nets for 3D data resort to either volumetric grids or collections of images.

## Challenges

This paper addresses the problem of 3D reconstruction from a single image, generating a straight-forward form of output – point cloud coordinates. Along with this problem arises a unique and interesting issue, that the groundtruth shape for an input image may be ambiguous. As Neural Network has to hallucinate complete 3D structure from a single view, the problem is under-determined hence opening the possibility of multiple possible outputs for a single 2D image. The problem is quite different from general regression problem where there is proper annotation for the training data and the output is deterministic. Here the ground truth is one possible 3D output out of all others equally probable on which net can be trained.

Due to ground truth ambiguity, the network design, the choice of loss function and the learning paradigm all become very crucial and their choice becomes highly sensitive in deciding the correctness of the output.

To deal with the above situation, the authors have chosen the following:

- A Conditional Generative Network
- Different Versions of Network Design
- Loss function of Chamfer Distance (Other than usual norms)
- Introduction of Random Vector 'r' & MoN Loss function as a wrapping function to deal with multiple outputs.

The reason for the above choices and their roles in getting the optimal output has been explained in the upcoming sections.

# Network Architecture

## Why Conditional Generative Network?

To deal with the challenges arose due to the choice of unorthodox output, multiple steps were taken by the authors, out of them one was to choose the architecture of the neural network to be Conditional Generative Network rather than the Discriminative Network. It is because this is a situation where there are multiple, equally good 3D reconstructions of a 2D image, making our problem very different from classical regression/classification settings, where each training sample has a unique ground truth annotation.

The network being generative tries to learn the whole distribution of the input image space and provides the output from the same distribution. The network is conditional sampler so that it samples plausible 3D point clouds from the estimated ground truth space given an input image.

Secondly, authors allow multiple reconstruction candidates for a single input image. This design reflects the fact that a single image cannot fully determine the reconstruction of a 3D shape.

## Network Formulation:

Goal - Single 2D image (RGB or RGBD) to complete 3D shape

Output Format -  $S = \{(x_i, y_i, z_i)\}_{i=1}^N$  ; where N is predefined constant,  $N = 1024$

Ground Truth - The ground truth is a probability distribution,  $P(.|I)$  over the shapes possible for input, I i.e. for a given image the output 3D points come with a probability. So, Neural Network is trained as a conditional sampler from  $P(.|I)$

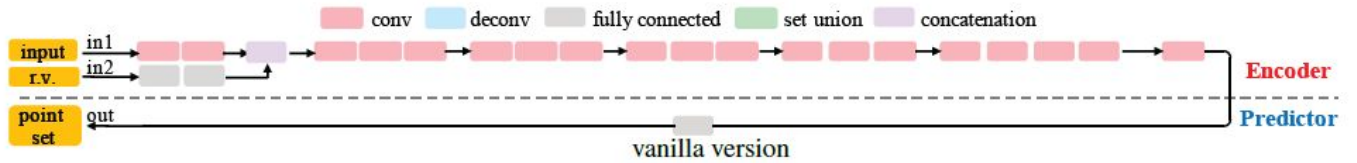
$$S = G(I, r, \theta)$$

$\theta$  - n/w parameter

$r \sim N(0, I)$  - It is a random variable to perturb the input

## Basic structure of the network:

### Vanilla Version:



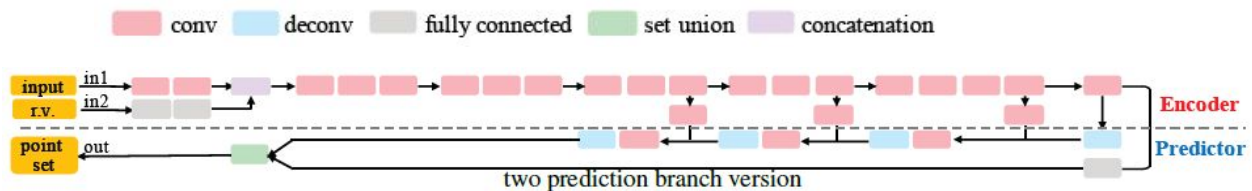
Network has an Encoder Stage and a Predictor Stage.

Encoder Stage: The encoder maps the input pair of an image  $I$  and a random vector  $r$  into an embedding space which can be understood by the predictor. The encoder is a composition of convolution and ReLU layers, in addition, a random vector  $r$  is subsumed so that it perturbs the prediction from the image  $I$ .

Predictor Stage: The predictor outputs a shape as an  $N \times 3$  matrix  $M$ , each row containing the coordinates of one point. It contains one fully connected layer.

Though simple, this version works reasonably well in practice. To further improve the design of the predictor branch to better accommodate large and smooth surfaces which are common in natural objects, a Two prediction branch version was designed.

### Two prediction branch version:



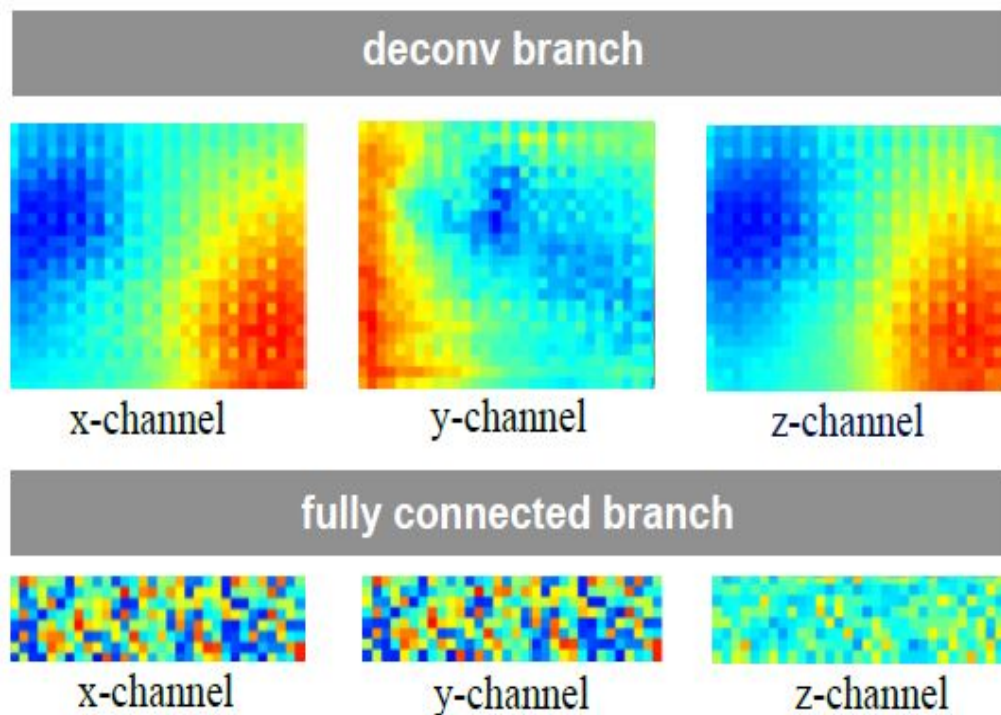
This version has two parallel predictor branches – a fully-connected (fc) branch and a deconvolution (deconv) branch. The fc branch predicts  $N1$  points as before. The deconv branch predicts a 3 channel image of size  $H \times W$ , of which the three values at each pixel are the coordinates of a point, giving another  $H \times W$  points. Their predictions are later merged together to form the whole set of points in  $M$ . Multiple skip links are added to boost information flow across encoder and predictor.

The fully connected predictor as above cannot make full use of such natural geometric statistics, since each point is predicted independently. The (conv-deconv) predictor exploits this geometric smoothness property. Due to which one enjoys high flexibility in capturing complicated structures(fully connected)and the other exploits geometric continuity(conv-deconv).

In summary, with the fc branch, our model enjoys high flexibility,showing good performance at describing intricate structures.

With the deconvolution branch, our model becomes not only more parameter parsimonious by weight sharing;but also more friendly to large smooth surfaces, due to the spatial continuity induced by deconv and conv.

### **Comparison of fully connected branch and deconv branch:**



In above Fig the values in the x, y and z channels are plotted as 2D images for one of the models. In the deconv branch the network learns to use the convolution structure to construct a 2D surface that wraps around the object. In the fully connected branch the output is less organized as the channels are not ordered. The deconv branch is in general good at capturing the “main body” of the object, while the fully connected branch complements the shape with more detailed

components. This reveals the complementarity of the two branches. The predefined weights sharing and node connectivity endow the deconv branch with higher efficiency when they are congruent with the desired output's structure. The fully connected branch is more flexible but the independent control of each point consumes more network capacity.

## Loss Function

A critical challenge is to design a good loss function for comparing the predicted point cloud and the groundtruth.

To plug in a neural network, a suitable distance must satisfy at least three conditions:

- 1) Differentiable with respect to point locations.
- 2) Efficient to compute, as data will be forwarded and back-propagated multiple times.
- 3) Robust against small number of outlier points in the sets.

We seek for a distance  $d$  between subsets in  $\mathbb{R}^3$ , so that the loss function takes the form:

$$L(\{S_i^{pred}\}, \{S_i^{gt}\}) = \sum d(S_i^{pred}, S_i^{gt}),$$

### Loss Function Definition:

**Chamfer distance** We define the Chamfer distance between  $S_1, S_2 \subseteq \mathbb{R}^3$  as:

$$d_{CD}(S_1, S_2) = \sum_{x \in S_1} \min_{y \in S_2} \|x - y\|_2^2 + \sum_{y \in S_2} \min_{x \in S_1} \|x - y\|_2^2$$

For each point, the algorithm of Chamfer Distance finds the nearest neighbor in the other set and sums the squared distances up. Viewed as a function of point locations in  $S_1$  and  $S_2$ , Chamfer Distance is continuous and piecewise smooth. The range search for each point is independent, thus trivially parallelizable. Also, spatial data structures like KD-tree can be used to accelerate nearest neighbor search. Though simple, Chamfer Distance produces reasonably high quality results in practice.



### Need of random vector 'r':

Facing the inherent inability to resolve the shape precisely, neural networks tend to predict a “mean” shape averaging out the space of uncertainty.

As shown in the Paper that mean shape kills the purpose of predicting the multiple shapes and it varies with the choice of the distance.

The ambiguity of groundtruth shape may significantly affect the trained predictor, as the loss function induces our model to predict mean of the possible shapes.

To better model the uncertainty or inherent ambiguity (e.g. unseen parts in the single view), we enable the system to generate distributional output.

Naively plugging  $G(I, r, \theta)$  into Loss Function to predict  $S_{\text{predi}}$  won't work, as the loss minimization will nullify the randomness. It is also unclear how to make CGAN work in our scenario, as building a discriminator that directly consumes a point set is itself an open problem.

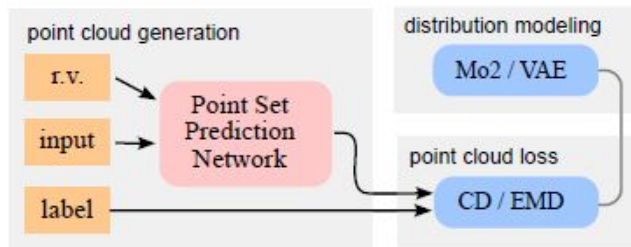
The problem can be solved by more complex frameworks like VAE, where we can incorporate secondary input channels (e.g. another view). However, we find practically a simple and effective method for uncertainty modeling:

The MoN loss: We train our network by minimizing a loss function as below:

$$\underset{\Theta}{\text{minimize}} \quad \sum_k \min_{\substack{r_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ 1 \leq j \leq n}} \{d(\mathbb{G}(I_k, r_j; \Theta), S_k^{\text{gt}})\}$$

Given an image  $I_k$ ,  $G(I, r, \theta)$  makes  $n$  predictions by perturbing the input with  $n$  random vectors  $r_j$ . Intuitively, we expect that one of the predictions will be close to the ground truth  $S_k^{\text{gt}}$  given by the training data, meaning that the minimum of the  $n$  distances between each prediction and the ground truth must be small.

We name this loss as Min-of-N loss (MoN), since it comes from the minimum of  $n$  distances. Any of the point set regression networks can be plugged into the meta network incorporating the MoN loss. In practice, we find that setting  $n = 2$  already enables our method to well explore the groundtruth space.



System Architecture

# Implementation Details

## How the dataset is prepared by them:

They rendered 2D views from CAD object models. Models are from the ShapeNet dataset, models with textures, containing large volumes of manually cleaned 3D objects. Concretely they used a subset of 220K models covering 2,000 object categories.

For each model, they normalized the radius of its bounding hemi-sphere to unit 1 and aligned their ground plane. Then each model was rendered into 2D images according to the Blinn-Phong shading formula with randomly chosen environmental maps. They used a simple local lightening model for the sake of computation time. However, it is straight-forward to extend our method to incorporate global illumination algorithms and more complex backgrounds.

## Implementation Details:

- Our network works on input images of 192x256.
- The deconv branch produces 768 points, which correspond to a 32x24 three channel image. The fully connected branch produces 256 points. So the final image is 32x32x3, where each location has 3 values depicting x,y,z.
- The training program is implemented in TensorFlow.
- 300000 gradient steps are taken, each computed from a minibatch of 32.
- Adam is used as the optimizer.
- We observed that the training procedure is smooth even without batch normalization.
- All activation functions are ReLU.

# Code and Results

## Network Architecture Code:

```
# Encoder branch
# first two conv layers
x=tflearn.layers.conv.conv_2d(x,16,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,16,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
# three conv layers
x=tflearn.layers.conv.conv_2d(x,32,(3,3),strides=2,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,32,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,32,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
# three conv layers
x=tflearn.layers.conv.conv_2d(x,64,(3,3),strides=2,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,64,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,64,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
# three conv layers
x=tflearn.layers.conv.conv_2d(x,128,(3,3),strides=2,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,128,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,128,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
# three conv layers
x=tflearn.layers.conv.conv_2d(x,256,(3,3),strides=2,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,256,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,256,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
# four conv layers
x=tflearn.layers.conv.conv_2d(x,512,(3,3),strides=2,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,512,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,512,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
x=tflearn.layers.conv.conv_2d(x,512,(3,3),strides=1,activation='relu',weight_decay=1e-5,regularizer='L2')
#last conv layer
x=tflearn.layers.conv.conv_2d(x,512,(5,5),strides=2,activation='relu',weight_decay=1e-5,regularizer='L2')
#Prediction branch
x_additional=tflearn.layers.core.fully_connected(x,2048,activation='relu',weight_decay=1e-3,regularizer='L2')
x_additional=tflearn.layers.core.fully_connected(x_additional,1024,activation='relu',weight_decay=1e-3,regularizer='L2')
x_additional=tflearn.layers.core.fully_connected(x_additional,256*3,activation='linear',weight_decay=1e-3,regularizer='L2')
x_additional=tf.reshape(x_additional,(BATCH_SIZE,256,3))
```

**Results:**

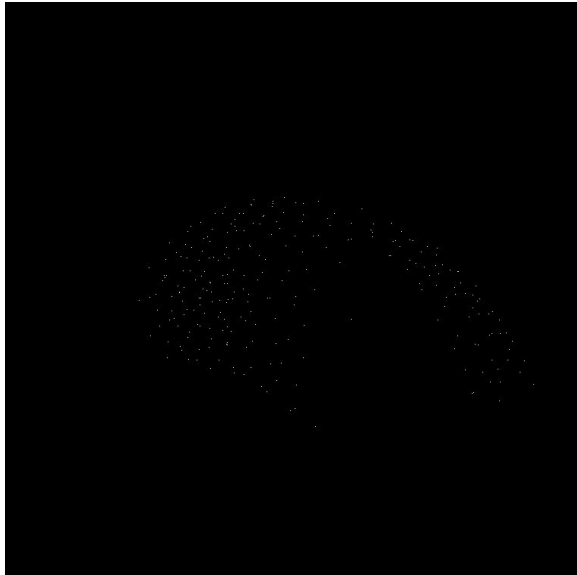
**Input**



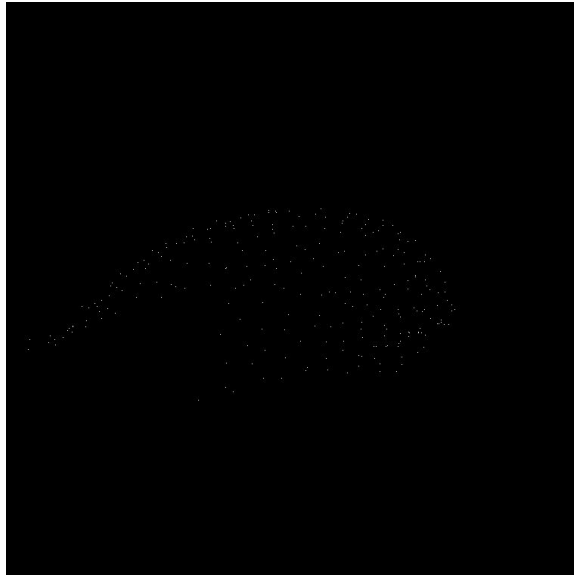
**Input\_Mask**



**Output View 1**



**Output\_View 2**



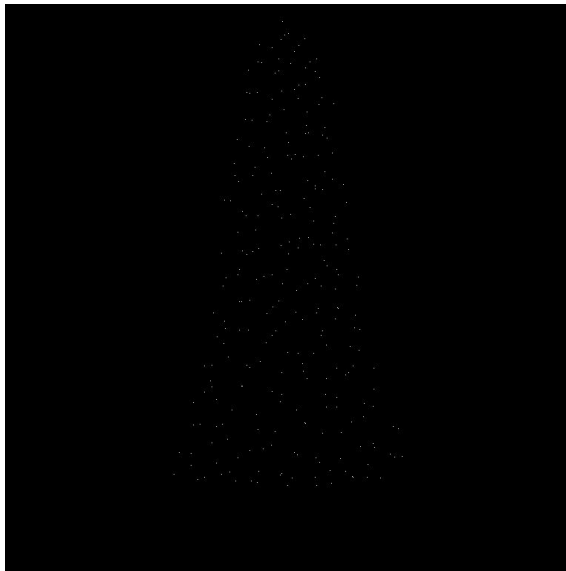
**Input**



**Input\_Mask**



**Output View 1**



**Output\_View 2**



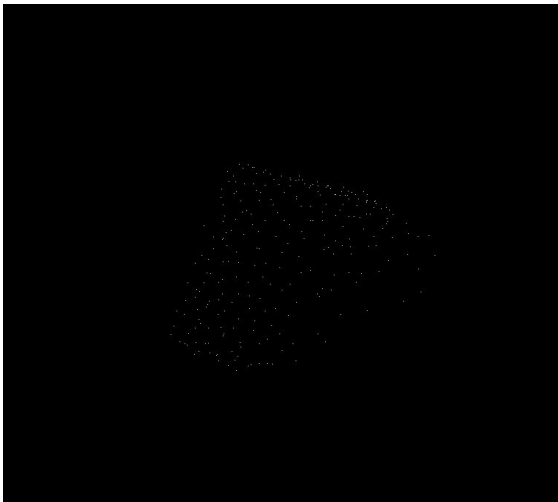
**Input**



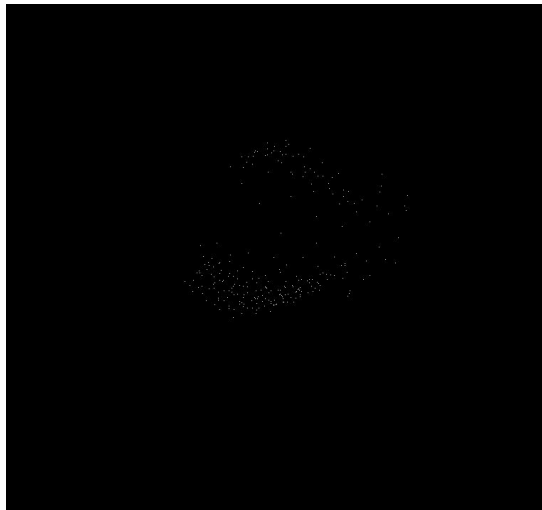
**Input\_Mask**



**Output View 1**



**Output\_View 2**



## **Future Tasks**

- Improvement of current results
- Implementation of Two Branch Prediction Version
- Comparison of results of Vanilla Version and Two Branch Prediction Version
- Testing the model on real world data

## References

- 1) [http://ai.stanford.edu/~haosu/papers/Sl2PC\\_arxiv\\_submit.pdf](http://ai.stanford.edu/~haosu/papers/Sl2PC_arxiv_submit.pdf)
- 2) M. Mirza and S. Osindero. Conditional generative adversarial nets. arXiv preprint arXiv:1411.1784, 2014.
- 3) ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], 2015
- 4) A unified approach for single and multiview 3d object reconstruction. arXiv preprint arXiv:1604.00449, 2016.
- 5) D. F. Fouhey, A. Gupta, and M. Hebert. Data-driven 3D primitives for single image understanding. In ICCV, 2013.