

Point Completion Networks and Segmentation of 3D Mesh

By

Naga Durga Harish Kanamarlapudi

February 2019

A Thesis/Project Submitted
in Partial Fulfillment
of the Requirements for the Degree of
Master of Science
in
Computer Engineering
Rochester Institute of Technology

Committee Approval:

Dr. Raymond Ptucha, *Advisor*

Associate Professor

Date

Dr. Alexander Loui, *Committee Member*

Professor

Date

Dr. Guoyu Lu, *Committee Member*

Associate Professor

Date

Acknowledgments

I would like to take this opportunity to thank my advisor Dr. Raymond Ptucha for his continual support and guidance during my master's degree. I am thankful for Dr. Alexander Loui and Dr. Guoyu Lu for being in my thesis committee. I would like to thank my family and my close friends, without their support this journey would not have been as joyful as it was.

Abstract

Deep learning has made many advancements in fields such as computer vision, natural language processing and speech processing. In autonomous driving, deep learning has made great improvements pertaining to the tasks of lane detection, steering estimation, throttle control, depth estimation, 2D and 3D object detection, object segmentation and object tracking. Understanding the 3D world is necessary for safe end-to-end self-driving. 3D point clouds provide rich 3D information, but processing point clouds is difficult since point clouds are irregular and unordered. Neural point processing methods like GraphCNN and PointNet operate on individual points for accurate classification and segmentation results. Occlusion of these 3D point clouds remains a major problem for autonomous driving. To process occluded point clouds, this research explores deep learning models to fill in missing points from partial point clouds. Specifically, we introduce improvements to methods called deep multistage point completion networks. We propose novel encoder and decoder architectures for efficiently processing partial point clouds as input and outputting complete point clouds. Results will be demonstrated on ShapeNet dataset.

Deep learning has made significant advancements in the field of robotics. For a robot gripper such as a suction cup to hold an object firmly, the robot needs to determine which portions of an object, or specifically which surfaces of the object should be used to mount the suction cup. Since 3D objects can be represented in many forms for computational purposes, a proper representation of 3D objects is necessary to tackle this problem. Formulating this problem using deep learning problem provides dataset challenges. In this work we will show representing 3D objects in the form of 3D mesh is effective for the problem of a robot gripper. We will perform research on the proper way for dataset creation and performance evaluation.

Contents

List of Figures	6
List of Tables.....	8
Chapter 1	9
1.1 Introduction	9
1.2 Motivation	10
1.3 Contributions	10
Chapter 2.....	11
2.1 Deep Learning	11
2.2 Convolutional Neural Networks	11
2.3 PointNet Related Architectures.....	14
A. PointNet.....	14
B. Dynamic Graph CNN.....	15
C. FoldingNet.....	16
D. 3D Point Capsule Networks.....	17
E. Point Completion Network.....	18
F. MeshNet.....	19
2.4 Self-supervised Learning.....	21
2.5 Metric Learning Loss Functions	23
A. Softmax and Cross Entropy Loss.....	23
B. Chamfer Distance Loss Function.....	23
C. Earth Mover Distance Loss Function.....	23
Chapter 3	24
3.1 Multi Stage Point Completion Network	224
A. Encoder	24
B. Voxel Feature Extraction Layer.....	24
C. Decoder	26
D. Training Strategy	27
3.2 Point Completion Network using Edge Convolution.....	28
A. Edge Convolution based Encoder.....	28

B.	Decoder	29
3.3	Capsule Based Point Completion Network	31
A.	Capsule Based Decoder.....	31
B.	Training Strategy	32
3.4	Multi-view Point Completion Network.....	33
A.	Mutli-view Encoder	33
B.	Decoder	34
C.	Training Strategy	35
3.5	Self-supervised Point Completion Network	35
A.	Encoder	35
B.	Decoder	36
C.	Training Strategy	37
3.6	Mesh Segmentation.....	35
A.	PointNet Center	37
B.	PointNet Mesh	38
C.	DGCNN Center	38
D.	MeshNet	39
	Chapter 4.....	40
4.1	Datasets	40
A.	ShapeNet	40
B.	Point Completion Networks	41
C.	Phenix Automation	42
	Chapter 5	46
5.1	Results	46
A.	MS PCN:	46
B.	Edge Convolution based PCN:	46
C.	Capsule-PCN:	47
D.	Multi-view PCN:.....	47
E.	Self-supervised PCN:	48
F.	PointNet Center:	48
G.	PointNet Mesh	50
H.	DGCNN Center	50
I.	MeshNet	50
6.1	Conclusions	53

List of Figures

Figure 1 An example of 1D Convolution [40].....	11
Figure 2 An example of 2D Convolution [41].....	12
Figure 3 An example of 2D Maxpooling [42].....	12
Figure 4 An example of 2D Average Pooling [43].....	13
Figure 5 Different Activation Functions [44].....	13
Figure 6 An example of Fully Connected Layer [45].....	13
Figure 7 PointNet architecture [15].....	14
Figure 8 DGCNN architecture [16].....	15
Figure 9 DGCNN spatial transform [16].....	15
Figure 10 DGCNN EdgeConv [16].....	16
Figure 11 FoldingNet architecture [34].....	17
Figure 10 Architecture of 3D Point Capsule Networks [46].....	18
Figure 11 Dynamic Routing algorithm [47].....	18
Figure 12 Architecture of PCN [2].....	19
Figure 13 MeshNet architecture [48].....	20
Figure 14 Face rotate convolution block [48].....	20
Figure 15 Mesh convolution block [48].....	21
Figure 16 An example of pretext task in self-supervised learning [35].....	22
Figure 17 Encoder for multistage point completion network.....	24
Figure 18 Voxel Feature Extraction Layer.....	25
Figure 19 Decoder for multistage point completion network	26
Figure 20 Stage-wise Learning.....	27
Figure 21 Stage-wise Learning.....	28
Figure 22 Stage-wise Learning.....	28
Figure 23 Stage-wise Learning.....	29
Figure 24 Capsule based decoder architecture.....	31

Figure 25 Multi-view encoder.....	33
Figure 26 Decoder for Multiview point completion network.....	34
Figure 27 Self-supervised pretext task for point completion network.....	35
Figure 28 Encoder for self-supervised learning.....	35
Figure 29 Decoder for self-supervised learning.....	36
Figure 30 PointNet Center.....	37
Figure 31 PointNet Mesh.....	38
Figure 32 DGCNN Center.....	38
Figure 33 MeshNet for segmentation.....	39
Figure 34 Different categories for ShapeNet dataset [8]	37
Figure 35 Examples of ShapeNet dataset [8]	40
Figure 36 Examples of PCN dataset [2]	41
Figure 37 Examples of Phenix Automation dataset	42
Figure 38 View of MATLAB tool for dataset annotation.....	43
Figure 39 Examples of annotated dataset for robot gripper.....	44
Figure 40 Results of PointNet-Center.....	47
Figure 41 Results of PointNet-Mesh.....	48
Figure 42 Results of DGCNN-Center.....	49
Figure 43 Results of MeshNet.....	50

List of Tables

Table 1 Results of MS PCN.....	46
Table 2 Results of Edge convolution based PCN.....	47
Table 3 Results of Capsule based decoder.....	47
Table 4 Results of Multiview DGCNN PCN.....	47
Table 5 Results of Self-supervised PCN.....	48

Chapter 1

Introduction

1.1 Introduction

Neural networks helped to obtain state of the art results in many of the challenging tasks like image recognition, image classification, 3D point cloud classification, text understanding and speech recognition. Convolutional Neural Networks (CNNs) have replaced the traditional computer vision techniques like Histogram of Oriented Gradients (HOG) and Scale Invariant Feature Transform (SIFT) to perform many of the vision related tasks like classification, and recognition. Architectures like ResNet [33], VGGNet [56], DenseNet [57] have shown the ability of CNNs to extract the features to obtain state of the art results in many vision related tasks like classification and segmentation.

Understanding the 3D world is one of the most challenging tasks in Artificial Intelligence (AI). Robust scene understanding is required in the case of applications like end-to-end autonomous driving and robotics. 3D data is typically represented in the format of point clouds and meshes. Due to the irregularity and unordered nature of point clouds, applying CNNs is not straight forward. Many of the previous works [1, 25, 13, 4, 17, 6] used hand crafted features of point clouds, other works [36] converted point clouds to voxels and used 3D CNNs to perform tasks, some works [10, 14] used Multiview CNNs by converting point clouds to images, and spectral CNNs [11, 12] use convolutions in spectral domain. Neural point processing using PointNet [15] extracts point cloud features by operating on individual point. PointNet++ [14] improves PointNet by extracting point cloud features hierarchically using multi-scale and multi-resolution grouping. Due to its simplistic architecture and effective performance, PointNet inspired many point cloud processing techniques [16, 5, 18].

Occlusion is a major problem in the real-world LiDAR scans. For safe end-to-end self-driving, incomplete point clouds should be made complete. PCN [2] takes the partial point cloud as input and outputs the complete point cloud. In this work, we propose novel encoder and decoder architectures for deep multistage point completion networks.

For a robot gripper like suction cup, in order to hold a 3D object firmly there should be a minimum flat surface area on the 3D object for the suction cup to get a good grip. The good area can be anywhere on the 3D object. Identifying different good parts on the outer surface of the 3D

object is difficult. 3D objects can be represented in different forms like point clouds, voxels and mesh. Representing 3D object with point clouds gives information of outer surface, but only in a sparse fashion. Representing 3D objects with voxels is generally limited by memory and doesn't have high resolution representation of the object surface. Further, the voxels in the center of the object are generally wasted memory. Representing 3D objects in the form of surface mesh provides more information of outer surface of 3D object and we can represent mesh in dense form due to its connectivity among different faces. So, formulating this problem as mesh segmentation helps to identify different parts of 3D object which are good a good surface for a robot gripper. Deep learning on 3D mesh data is a newer problem, [48] is the first deep learning model applied to 3D mesh data. After formulating this problem as a machine learning problem, getting the dataset suitable for this task is very difficult because there are no publicly available datasets dealing this problem. In this work, we show some methods for creating and annotating the dataset effectively and show initial results using existing deep learning architectures.

1.2 Motivation

Understanding the 3D world is very important for end-to-end autonomous driving, and occlusions create many problems as deep learning networks typically only see only a partial shape of the object. This research completes the partial shape of the object to minimize this problem. Also, for efficient point cloud registration, registration on full shapes provides better registration results instead of partial shapes. This work deals with completing a partial point cloud.

To identify good and bad faces of 3D objects to hold firmly by a robot gripper such as a suction cup, there are no existing algorithms. Converting the 3D object into 3D mesh and segmenting the mesh into good and bad faces provides information regarding various good faces for a robot gripper like suction cup to hold firmly.

1.3 Contributions

The main contributions of this thesis work can be summarized as follows:

- Experimented with different novel deep learning architectures.
- Combining the voxel wise and point wise features for better global feature vector.
- Complete shape of the input partial point cloud is optimized in multiple stages like coarse, middle and fine.
- Explored the capsule based dynamic routing architecture for the problem of point completion network.

- Experimented with self-supervised approach for the task of point completion network.
- Created MATLAB tool for dataset generation for the robot gripper task.
- Experimented with different deep learning architectures for the robot gripper task.

Chapter 2

Background

2.1 Deep Learning

Deep learning has achieved state of the art results in many computer vision related tasks like 2D and 3D object detection, mesh classification, point cloud classification and segmentation, and image captioning. Deep learning has replaced traditional computer vision techniques for feature extraction in the case of 3D data like point clouds and meshes. [36, 55] converts point clouds to voxels, [15, 16, 14] computes point-wise features and applies symmetric function like maxpooling, and [48] extracts the mesh features by considering each face as a unit.

2.2 Convolutional Neural Networks

Convolutional Neural Networks (CNNs) have become the standard method for extracting the features of images and point wise features in point clouds. CNNs typically consist of the following layers:

- Convolution layer
- Pooling layer
- Activation layer
- Fully connected layer

A convolution layer consists of multiple filters of a window size multiplied with the input pixels and performs a linear combination of all the multiplied values in that filter window. Filters in the initial layers learn low level features and filters in the later layers learn higher level features of the input data. Usually, the convolution operation will be performed in 1D, 2D or 3D. A 1D convolution is shown in Figure 1 and a 2D convolution operation is shown in Figure 2.

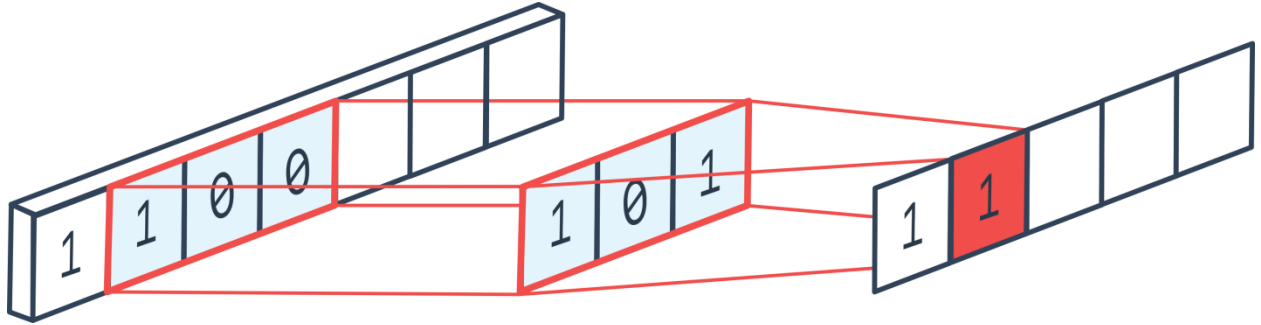


Figure 1 An example of 1D Convolution [40].

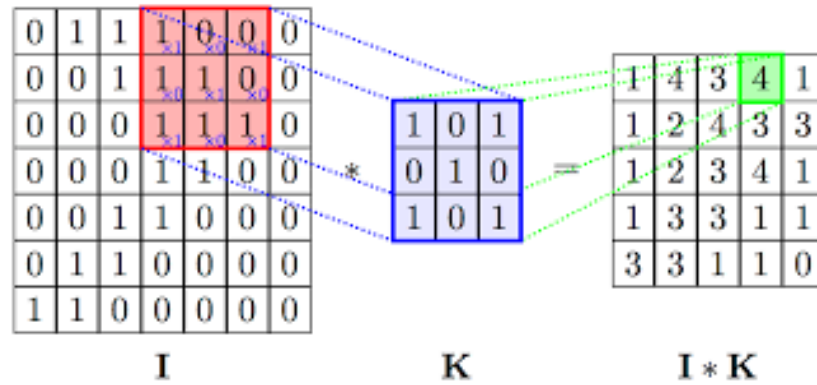


Figure 2 An example of 2D Convolution [41].

A pooling layer down samples the image depending on the type of pooling used. There are two kinds of pooling. One is max pooling as shown in Figure 3, which outputs the maximum value of the pixels in a certain pooling window. The other is average pooling as shown in Figure 4, which outputs the average value of all pixel values in a certain pooling window.

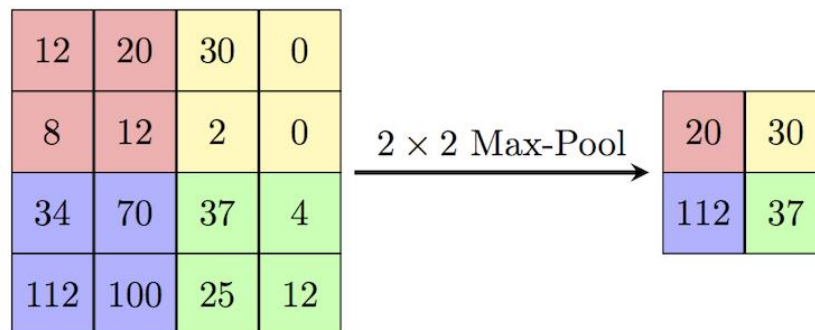


Figure 3 An example of 2D Maxpooling [42].

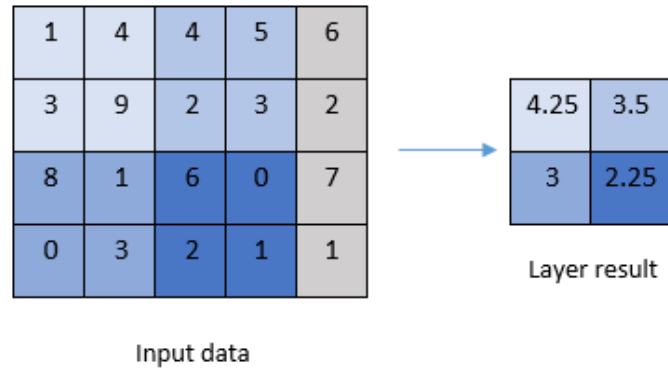
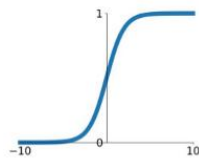


Figure 4 An example of 2D Average Pooling [43].

An activation layer performs a non-linear operation in a certain way depending on the type of activation function used. There are several activation functions like sigmoid, tanh, ReLU, and leaky ReLU. Figure 5 shows different activation functions used in deep learning.

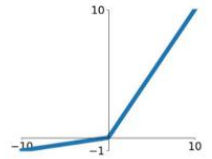
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



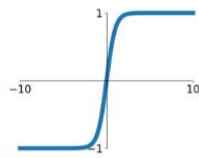
Leaky ReLU

$$\max(0.1x, x)$$



tanh

$$\tanh(x)$$

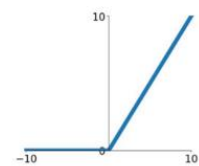


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ReLU

$$\max(0, x)$$



ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

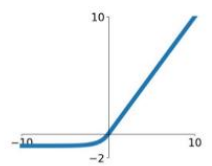


Figure 5 Different Activation Functions [44].

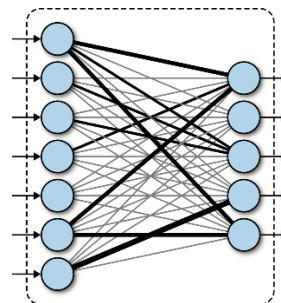


Figure 6 An example of Fully Connected Layer [45].

A fully connected layer is used to transform a high dimensional representation into a n -dimensional representation by connecting all the neurons in the input and fully connected layer. An example of fully connected layer is shown in Figure 6.

2.3 PointNet related architectures

A. PointNet:

PointNet [15] operates on the point clouds directly without converting them into other forms like voxels, a vector representation or rendering multiple views from point clouds. In order to account for the unordered nature of point clouds, PointNet proposed to use a symmetric function like maxpool to aggregate the features from the point clouds. PointNet consist of multi-layer perceptron (MLPs) to learn the point features and then apply symmetric function like maxpool to aggregate the point-wise features.

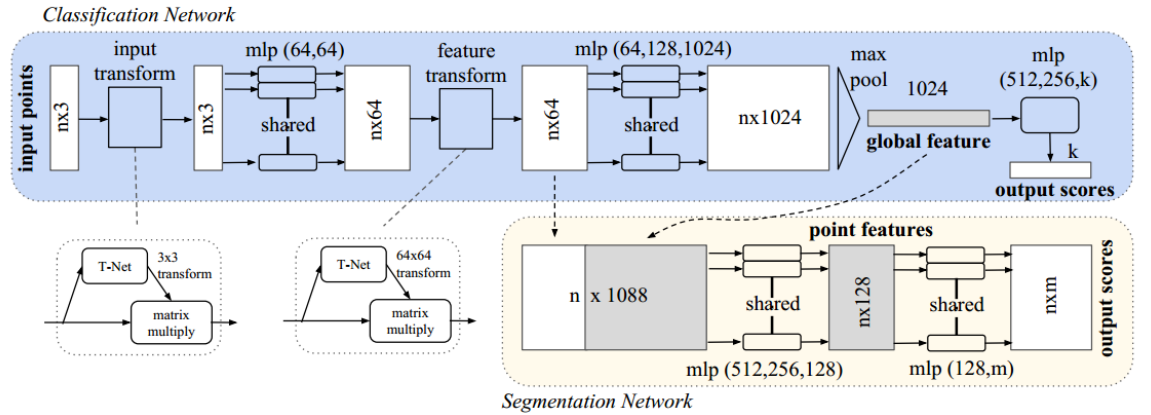


Figure 7 PointNet architecture [15].

PointNet consist of input transform and feature transform which are mini PointNet [15] like architectures. Input transform provides a 3×3 matrix which is applied on every point in the input point cloud and similarly feature transform provides a 64×64 matrix which is multiplied with pointwise features. A global feature vector is obtained by applying a symmetric maxpooling function such that this vector can be used in both a classification and segmentation network. The classification network outputs ' k ' score and the segmentation network outputs ' $n \times m$ ' scores. There are many PointNet inspired architectures. The architecture of [15] is shown in Figure 7.

B. Dynamic Graph CNN:

Inspired by [15], [16] also operates on individual points to learn salient features. [15] considers each point individually and does not consider the relationships between point pairs. To account for this DGCNN constructs a knn graph and learns the local properties between the point pairs. DGCNN applies multi-layer perceptron on the k nearest graph constructed and applies symmetric maxpooling to get the aggregated features of the graph.

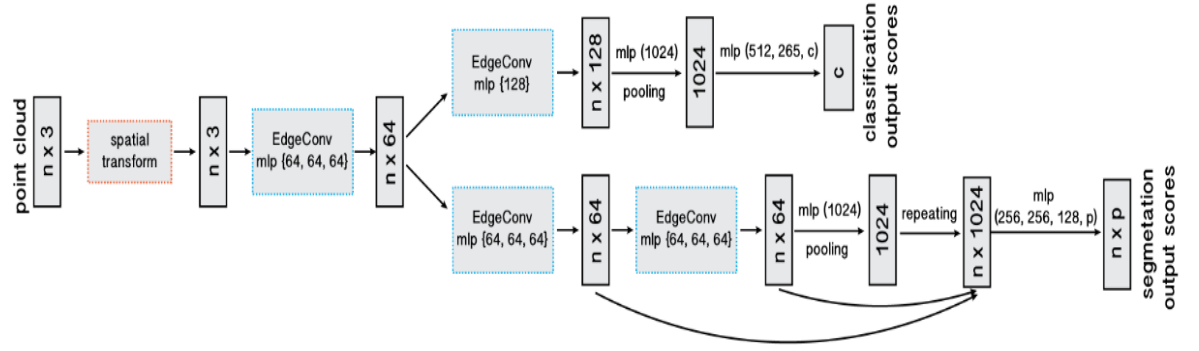


Figure 8 DGCNN architecture [16].

Figure 8 shows the DGCNN architecture, consisting of spatial transform and edge convolution block. The network consists of both a classification block and a segmentation block.

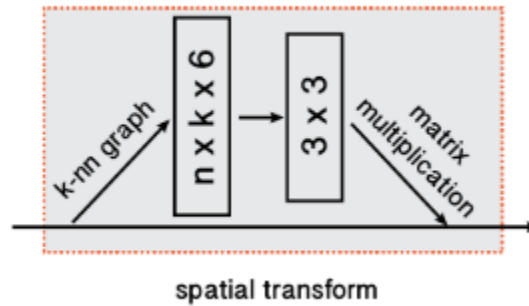


Figure 9 DGCNN spatial transform [16].

The spatial transform block aligns the input point cloud into a canonical space by applying a 3×3 matrix. This 3×3 matrix is estimated by constructing a k -nearest graph and then concatenating the point features with its k -nearest point features.

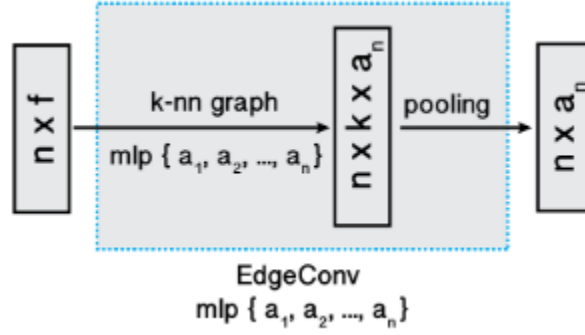


Figure 10 DGCNN EdgeConv [16].

The most important part of DGCNN is the edge convolution. Edge convolution is performed by constructing a k -nearest graph, then applying a multi-layer perceptron and then applying maxpooling to obtain the local features. The edge convolution operation is shown in (1).

$$e_{ijm} = \text{ReLU} \left(\Theta_m \cdot (x_j - x_i) \cdot (x_i) \right) \quad (1)$$

In (1), x_i is each point in the input point cloud, x_j is each point in the k -nearest graph of x_i and Θ_m is the weight matrix which can be approximated by a multi-layer perceptron.

C. FoldingNet:

FoldingNet is a point cloud auto-encoder. FoldingNet [34] consists of a graph-based encoder on top of [15] and a folding-based decoder. FoldingNet [34] deforms a 2D grid into a 3D object surface of the point cloud.

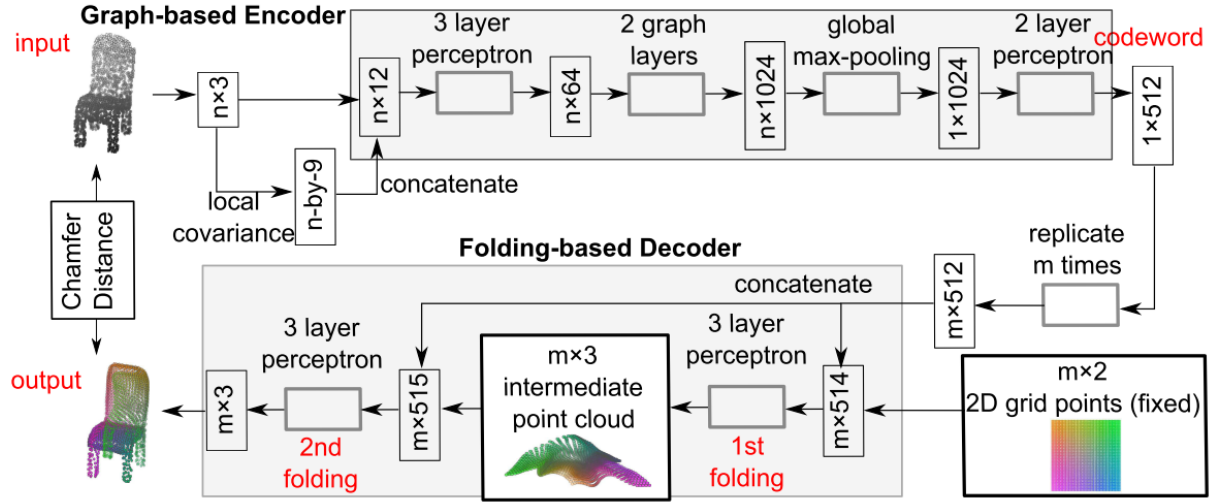


Figure 11 FoldingNet architecture [34].

The graph based encoder consists of multi-layer perceptron and graph based maxpooling layers. First a local covariance matrix of size 3×3 is computed using 3D positions of the points and its one hop neighbors. For every point in the input point cloud, the local covariance matrix is constructed to give a vector of $n \times 9$ and is concatenated with the input points of size $n \times 3$ to give a matrix of size $n \times 12$. The graph is a k -nearest graph computed by considering the k nearest neighbors of the 3D points and after that a maxpooling operation is performed to aggregate the features. A final codeword of size 1×512 is obtained as an output from the encoder.

The folding based decoder deforms the 2D fixed grid points to a 3D point cloud. The codeword from the graph-based encoder is fed into the decoder by replicating it m times and concatenated with the 2D grid points to obtain a matrix of size $m \times 514$. A multi-layer perceptron is applied on this matrix by processing it row-wise to obtain an intermediate point cloud. This concatenation and point-wise convolution is applied again on this intermediate point cloud to obtain a final output.

D. 3D Point Capsule Networks:

3D point capsule networks are an auto-encoder. They consist of dynamic routing between primary point capsules and latent capsules. Individual feature maps are computed by applying point-wise

multi-layer perceptron. Primary point capsules are obtained by feeding these feature maps into multiple convolutional layers with different weights. Latent capsules are obtained by applying a dynamic routing algorithm to the primary point capsules. These latent capsules attend to different parts of the input point clouds because of the dynamic routing algorithm.

The decoder concatenates a 2D random grid of fixed size to the latent capsules obtained from the encoder. In order to reconstruct the point cloud of input size, latent capsules are replicated m times and then concatenated with a 2D random grid.

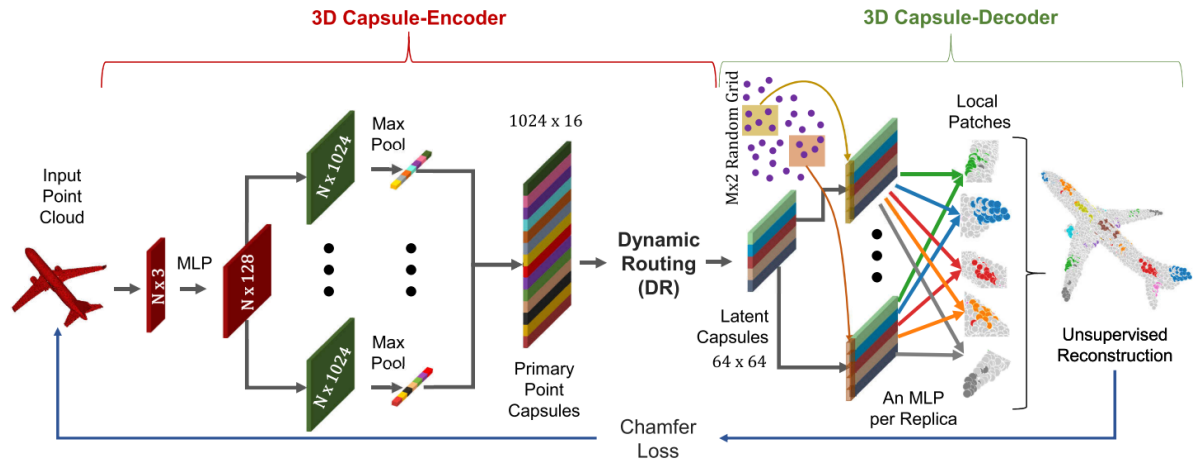


Figure 10 Architecture of 3D Point Capsule Networks [46].

Procedure 1 Routing algorithm.

```

1: procedure ROUTING( $\hat{\mathbf{u}}_{j|i}, r, l$ )
2:   for all capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $(l + 1)$ :  $b_{ij} \leftarrow 0$ .
3:   for  $r$  iterations do
4:     for all capsule  $i$  in layer  $l$ :  $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$ 
5:     for all capsule  $j$  in layer  $(l + 1)$ :  $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{\mathbf{u}}_{j|i}$ 
6:     for all capsule  $j$  in layer  $(l + 1)$ :  $\mathbf{v}_j \leftarrow \text{squash}(\mathbf{s}_j)$ 
7:     for all capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $(l + 1)$ :  $b_{ij} \leftarrow b_{ij} + \hat{\mathbf{u}}_{j|i} \cdot \mathbf{v}_j$ 
   return  $\mathbf{v}_j$ 

```

Figure 11 Dynamic Routing algorithm [47].

Dynamic routing algorithm is applied as shown in Figure 11. The algorithm is applied between all the point capsules in primary point capsules and all the capsules in latent capsules.

E. Point Completion Network:

Point completion network (PCN) is an encoder decoder architecture which completes an input partial point cloud and produces a complete point cloud. The encoder of [2] applies a point-wise multi-layer perceptron and then a symmetric function such as maxpooling to obtain a global feature vector. This global feature vector is replicated m times, concatenated with point-wise features and then applied multi-layer perceptron. A final global feature vector is computed by applying maxpooling. The final global vector is of size 1024.

The decoder of PCN [2] takes global feature vector from the encoder as input and computes a coarse output by applying a multi-layer perceptron. The coarse output is concatenated with a 2D grid of fixed size with radius 4 and with the m times global feature vector. The shared multi-layer perceptron is applied to output a final detailed output.

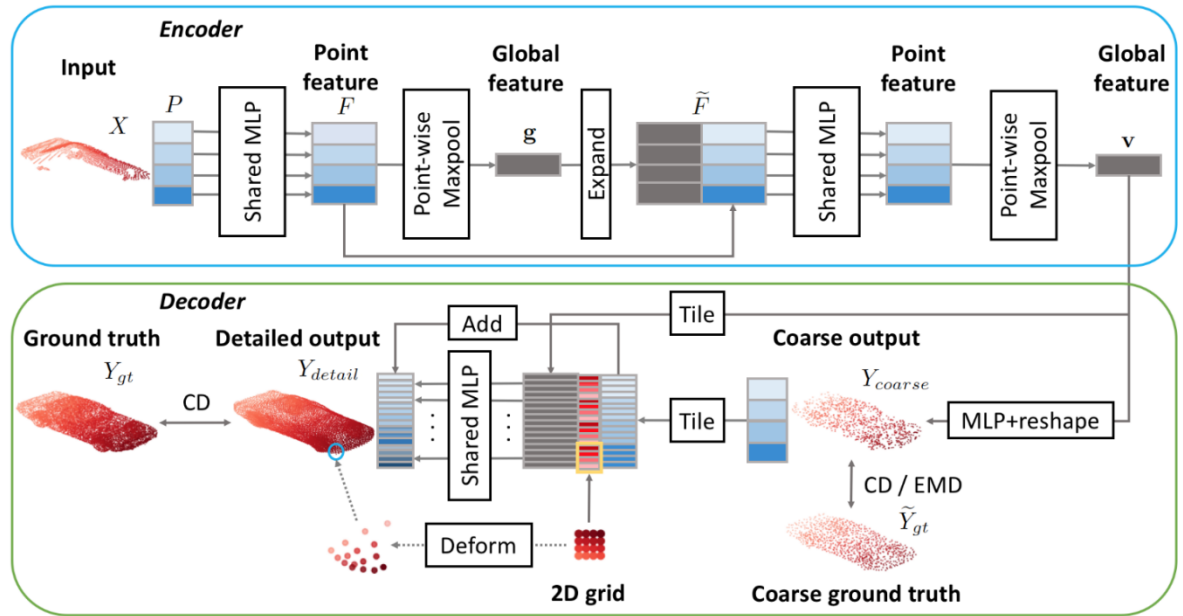


Figure 12 Architecture of PCN [2].

F. MeshNet:

MeshNet [48] is a neural network applied on 3D mesh data, considering each face of the 3D mesh as an operating unit. MeshNet [48] consist of a mesh convolution block, a structural descriptor and

a spatial descriptor. The spatial descriptor takes centers as inputs and then applies a multi-layer perceptron to compute the center features. The structural descriptor takes neighboring indices, the normal of each face and the corners of each face as inputs, then computes features using face kernel correlation and face rotate convolution.

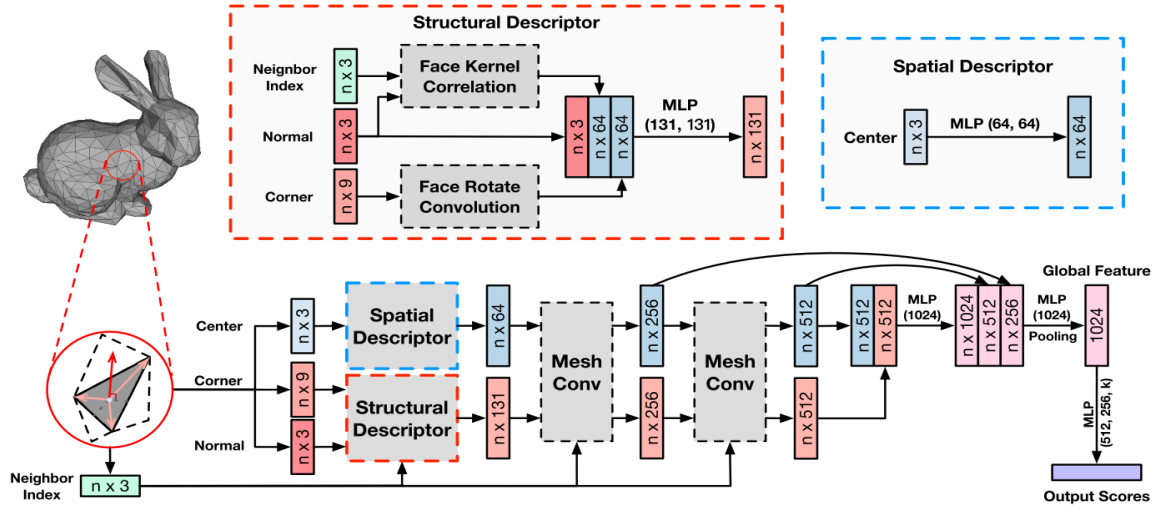


Figure 13 MeshNet architecture [48].

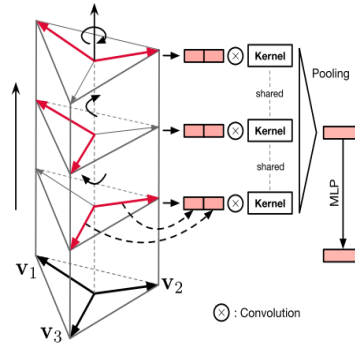


Figure 14 Face rotate convolution block [48].

The face rotate convolution block as shown in Figure 14 rotates the face and then applies a convolution operation on pairs of corner vectors. The face kernel correlation applies correlation between neighboring corners and a Gaussian kernel.

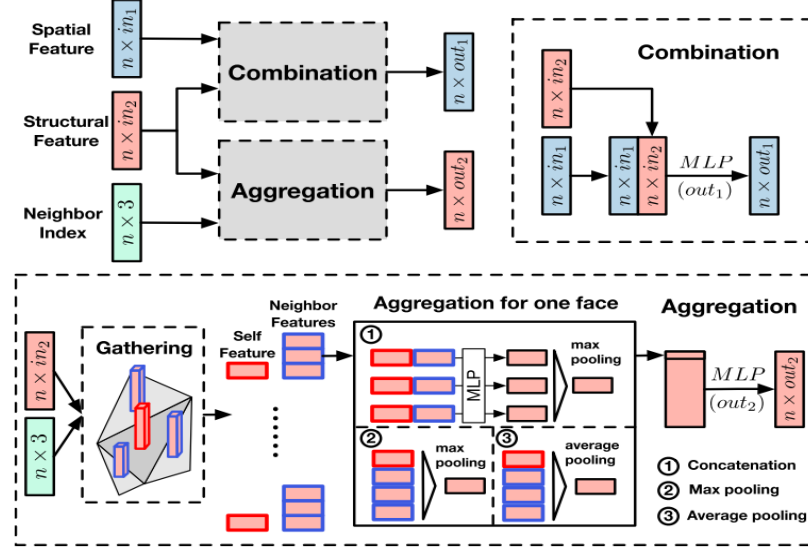


Figure 15 Mesh convolution block [48].

The mesh convolution block consists of a combination block and an aggregation block. The combination block combines the spatial features and structural features by concatenating spatial and structural features. The aggregation block takes neighboring indices and structural features as input, then computes the aggregated features by applying multi-layer perceptron and maxpooling operation.

2.4 Self-supervised Learning:

Transfer learning always helps the model if the task has less training data. But in many tasks, such as medical imaging and point completion, we cannot find a model to transfer the weights from. In tasks where transfer learning is not possible, self-supervised learning plays a key role in learning the input dataset features. In self-supervised learning we define a task known as a pretext task. We train our model and use these pretext task weights to perform transfer learning on our real, but limited training data. These pretext tasks don't need any ground truth data, rather the model defines a task and extracts its own ground truth for initial training of our model.

Different pretext tasks are defined in [58], [59], [60], [35]. [58] defined the task of colorization as the pretext task, [59] defined guessing the spatial information by randomly sampling the different patches of an image as the pretext task, [35] rotated the image into four different directions and defined classifying the angle rotation as the pretext task.

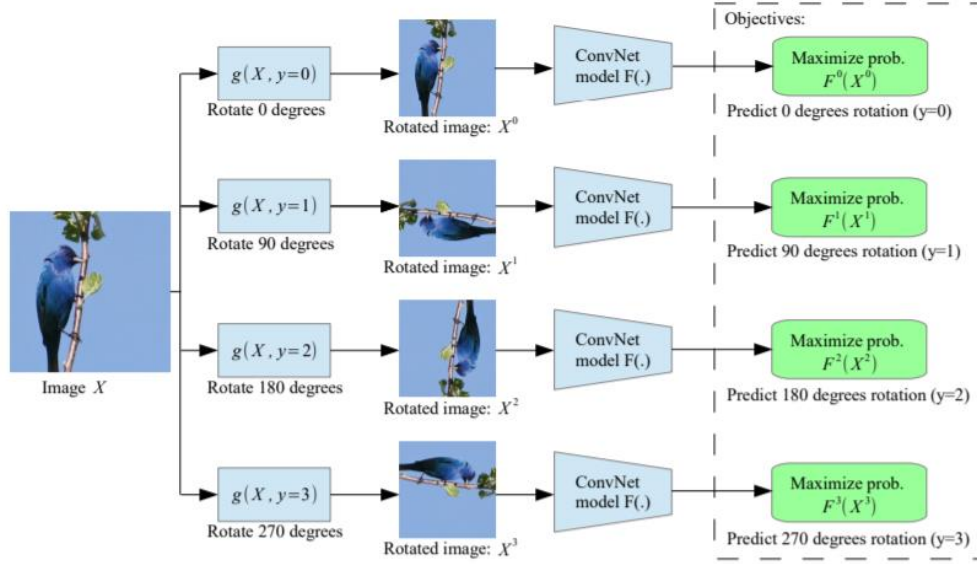


Figure 16 An example of pretext task in self-supervised learning [35].

The pretext task used in [35] is shown in Figure 16, which shows the pretext task of rotating the image by different angles and then predicting the rotated class. Simple tasks like this helps the model to learn semantic features in a robust way. After pretext task learning, we use these weights to perform transfer learning on our ground truth labelled data.

2.5 Loss Functions

A. Softmax and Cross Entropy Loss:

The softmax activation function takes an N -dimensional real vector as input and outputs an N -dimensional real vector with values in the range $(0, 1)$ that add up to 1. It is usually used as the last layer before calculating loss.

$$S_j = \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}} \quad \forall j \in 1 \dots N \quad (2)$$

Softmax uses (2) to calculate the probability of each class prediction in the input vector.

Cross entropy loss calculates the distance between the output of deep neural network and the original distribution. Cross entropy loss is calculated using (3).

$$H(y, p) = - \sum_i y_i \log p_i \quad (3)$$

B. Chamfer Distance Function:

$$\begin{aligned} CD(S_1, S_2) = & \frac{1}{|S_1|} + \sum_{x \in S_1} \min_{y \in S_2} ||x - y||_2 + \frac{1}{|S_2|} \\ & + \sum_{y \in S_2} \min_{x \in S_1} ||y - x||_2 \end{aligned} \quad (4)$$

The Chamfer Distance (CD) as shown in (4) calculates the average distance between the input point cloud and output point cloud. In (4) S_1 denotes the predicted complete point cloud and S_2 denotes the ground truth complete point cloud.

C. Earth Mover Distance (EMD) Loss Function:

$$EMD(S_1, S_2) = \min_{\phi: S_1 \rightarrow S_2} \frac{1}{|S_1|} + \sum_{x \in S_1} ||x - \phi(x)||_2 \quad (5)$$

The Earth Move Distance (EMD) as shown in (5) finds a bijection from the input point cloud to the output point cloud, minimizing the average distance between corresponding points. Here bijection function finds the minimum cost to move from the predicted point coordinates to the actual ground truth coordinates.

3.1 Multistage point completion network

A. Encoder:

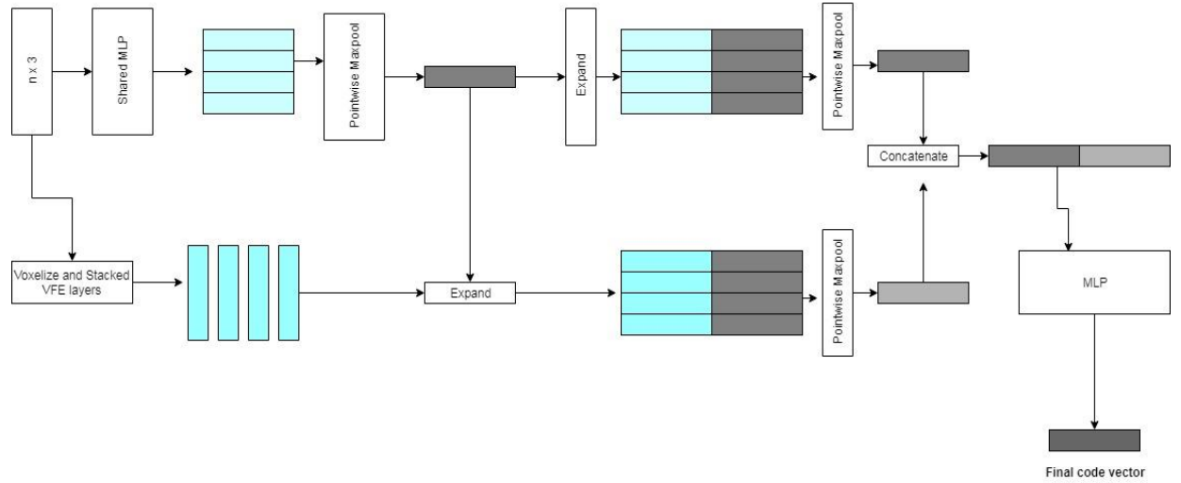


Figure 17 Encoder for multistage point completion network.

Figure 17 shows the proposed encoder for point cloud completion. The Point Completion Network (PCN) [2] extracts the point-wise features using style shared multi-layer perceptron, concatenates point-wise features with global features and then applies shared multi-layer perceptron layers. Many 3D object detection networks [36,55] voxelize the input point cloud and then performs multi-layer perceptron. Since our input is a partial point cloud, interacting different voxels and global features is desirable to extract the features in a better way. Input point clouds are voxelized and voxel features are extracted using voxel feature extraction (VFE) layers.

B. Voxel Feature Extraction Layer:

Voxelization is converting the input point cloud into a 3D grid, whereby each sub-grid is called a voxel. Each voxel is assigned fixed number of points belonging to that voxel. If a voxel doesn't contain minimum of number of points padding is used. VFE layers are used to extract the voxel-wise features.

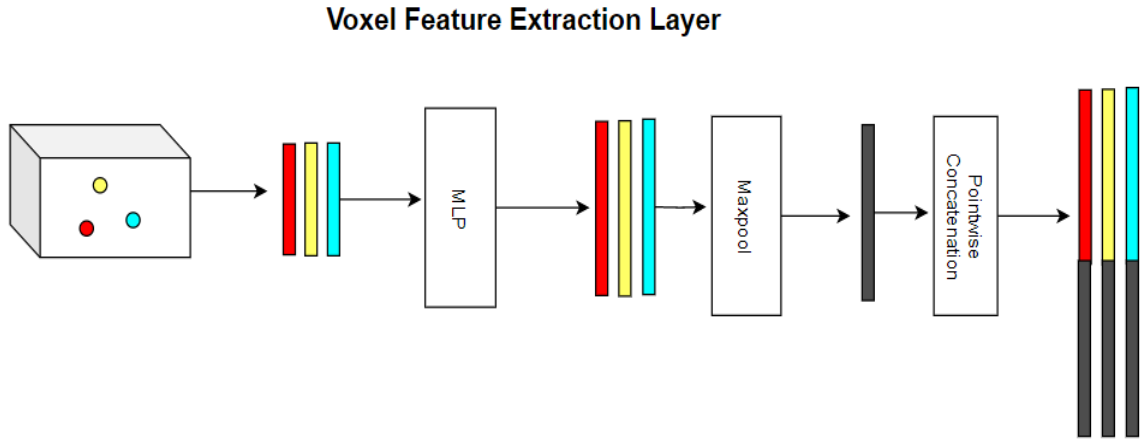


Figure 18 Voxel feature extraction Layer.

After extracting voxel-wise features, these voxel-wise features are concatenated with global features extracted by a shared multi-layer perceptron. This ensures learning the relationships between different points in the input point cloud. Without voxelization, it is difficult to learn the relationships in the input partial point cloud. After applying a shared multi-layer perceptron to the concatenated voxel feature vector and global feature vector, a symmetric maxpooling operation is applied. This vector is concatenated with the other global feature vector and then passed through a shared multi-layer perceptron. A final vector of size 1024 is obtained as an output from the encoder. An example of voxel feature extraction layer is shown in Figure 18.

C. Decoder:

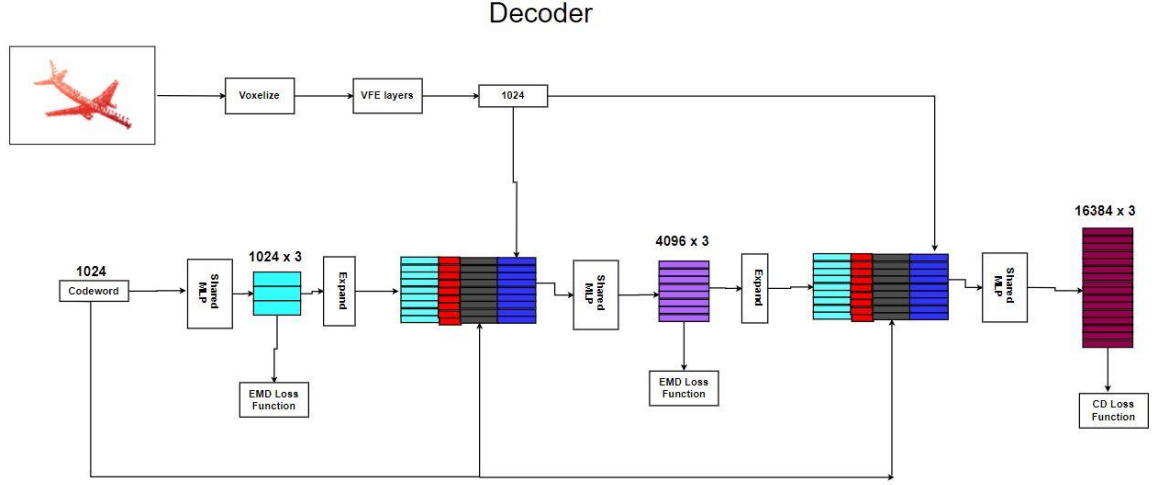


Figure 19 Decoder for multistage point completion network.

Our proposed decoder is shown in Figure 19. PCN [2] uses the codeword from the encoder as input and optimizes the codeword in the form of a coarse and fine output. The decoder shown in Figure 19 optimizes the code in 3 stages- in the form of a coarse output, middle output and final output. Since the problem is completing the input partial point cloud, a codeword from the encoder will have global information. To use this information effectively, the model needs to learn how to optimize the different sizes of point clouds. The proposed decoder outputs a coarse point cloud, a middle point cloud and a fine complete point cloud.

As shown in Figure 19, the shared multi-layer perceptron is applied to the codeword obtained from the encoder to output a coarse output of size 1024×3 . FoldingNet [34], AtlasNet [54] and PCN [2] deform a 2D grid to the 3D point cloud. Our decoder deforms the 2D grid in multiple stages. The coarse output is concatenated with a tiled fixed 2D grid and tiled codeword to produce a middle output of size 4096×3 . The middle output is concatenated with another tiled fixed 2D grid and a tiled codeword to output a final complete point cloud. The coarse output and middle output are optimized using the EMD (4) loss function, while the final output is optimized using the CD (3) loss function.

D. Training Strategy

The encoder and decoder are coded in the TensorFlow framework. The model is trained for 15 epochs with an initial learning rate of 0.0001 and uses exponential decay with decay rate 0.1. The learning rate is decayed at several steps during training. The model uses batch normalization [53] at each multi-layer perceptron layer. Training was performed in two different ways. In the first way the model is trained end to end, optimizing the losses of coarse output, middle output and complete point clouds output. During initial iterations of training the middle output and complete output losses are not as important as the initial coarse output. This is done using two parameters, α and β . Starting from value 0.01, both parameters values increased after every 'm' iterations.

In the second way of training, the model is divided into three stages. The stage 1 model is trained only for the coarse output. This coarse output is taken by a stage 2 model and trained for middle output. The stage 3 model takes middle output as input and trained only for complete point cloud. Figure 20 shows the stage-wise learning procedure.

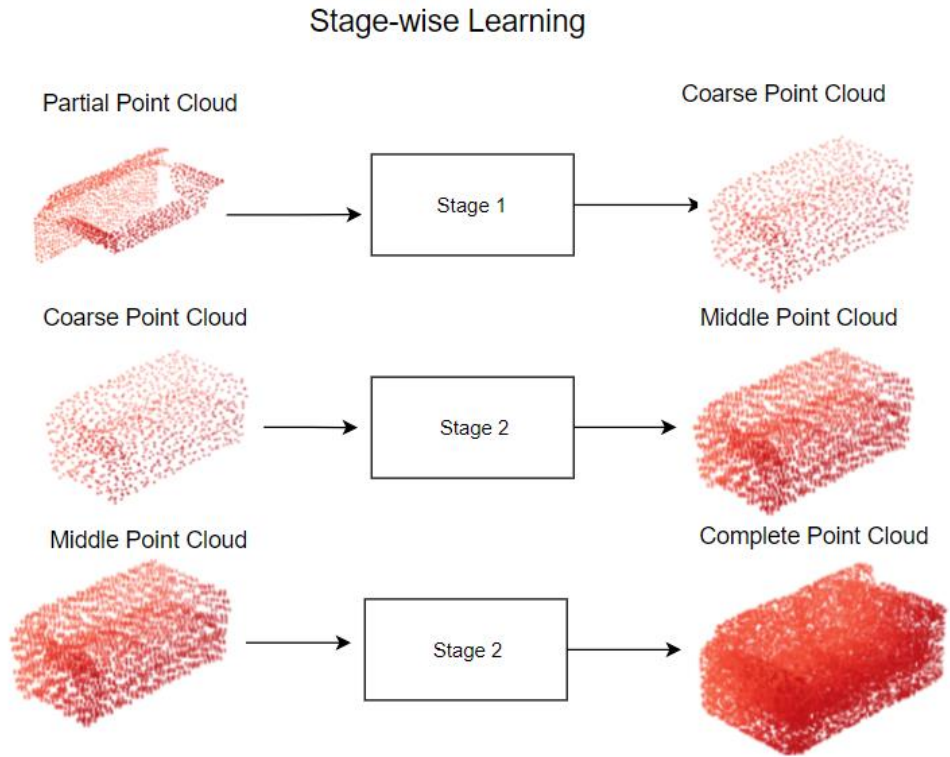


Figure 20 Stage-wise Learning.

3.2 Point Completion Network using Edge Convolution:

A. Edge convolution based network:

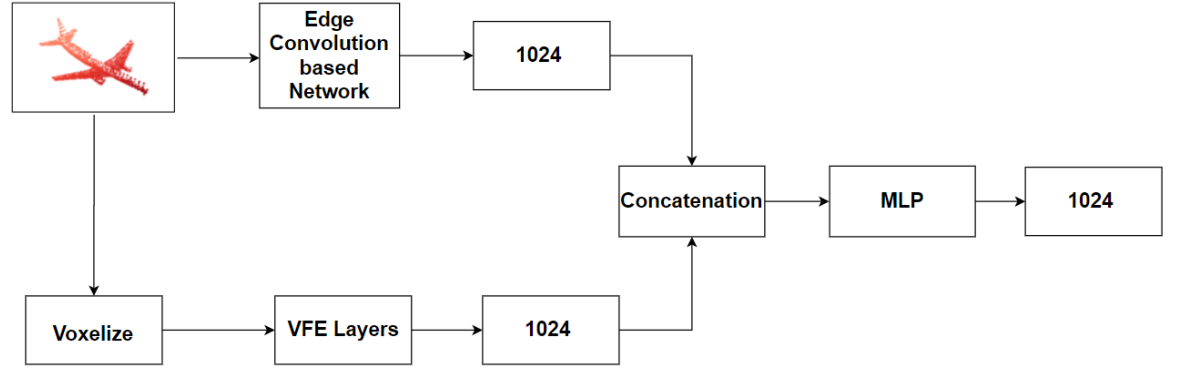


Figure 21 Edge convolution based encoder architecture.

DGCNN [16] employs edge convolution to effectively extract the local features by constructing a k -nearest graph, applies multilayer perceptron and then maxpooling operation. PCN [2] extracts the pointwise features by using PointNet [15] style architecture. The architecture of the edge convolution based point completion network is shown in Figure 21. The input partial point cloud is voxelized and extracted voxel features by using VFE (Voxel Feature Extraction) layers and then concatenated with the final vector obtained from the edge convolution network, and then applied multiplayer perceptron to obtain the final code vector.

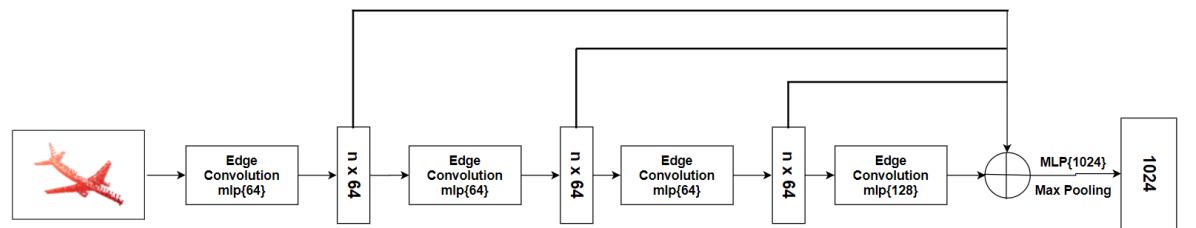


Figure 22 Complete architecture of edge convolution network.

The detailed architecture of the edge convolution network is shown in Figure 22. The architecture contains 4 edge convolution blocks, the output of the edge convolution blocks are concatenated, applied multilayer perceptron and then maxpooling to obtain the full 1024 vector. The input to the network is a partial point cloud of shape 2048×3 . The operation of the edge convolution is same as in [16]. The input to the VFE layers is of shape $m \times p \times 3$, where m is the number of voxels, p is the maximum number of points in the voxel. Interacting the points in different voxels is necessary to learn the information of missing points in the input point cloud. Because of this reason we employed 2D convolutions to extract the features from the voxels. The final vector from the VFE layers is of shape 1024.

B. Decoder:

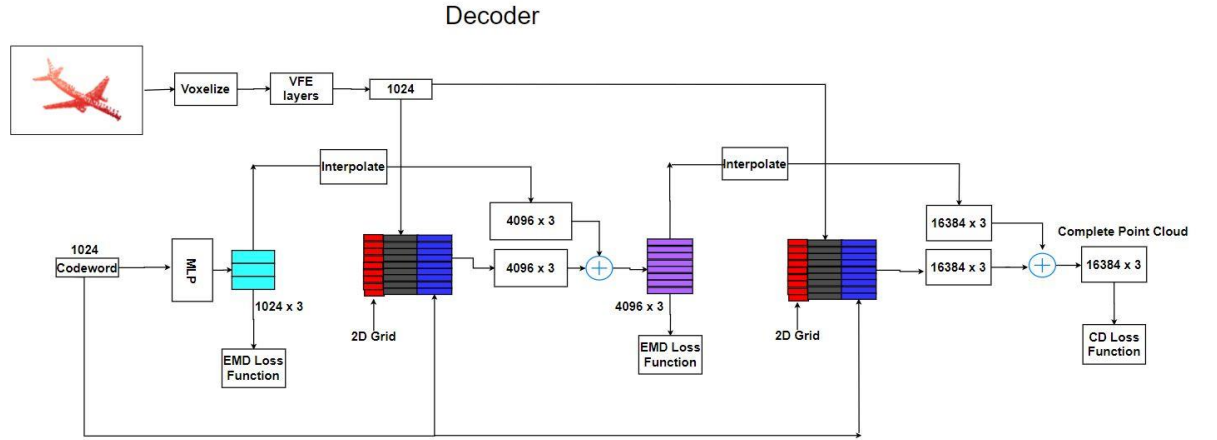


Figure 23 Architecture of the decoder employing interpolation.

The modified decoder which employs interpolation technique instead of tiling is shown in Figure 23. The decoder optimizes the complete point cloud in multiple stages: the coarse stage, the middle stage and the complete stage. The ground truths for each of the stages are generated by using FPS (Farthest Point Sampling) which is used in [14]. FPS algorithm is good compared to random sampling in sampling the point clouds. Both, the middle output and the complete point cloud are generated by concatenating a 2D grid of fixed size to the code word from the encoder and the final vector from the VFE layers. The coarse and the middle point are optimized using EMD loss function, and the complete point cloud is optimized using CD loss function.

In [2] concatenates the tiled coarse point cloud to the 2D random grid to generate the complete point cloud. A better way to use the information of the coarse point cloud is an interpolation, in the tiling, duplication of the same features is happening which doesn't help that much for the generation of the complete point cloud. The decoder used in this architecture interpolates the coarse point cloud to the shape of the middle point cloud and added to the middle point cloud features generated using the folding technique. Similarly, the middle point cloud is interpolated to the shape of the complete point cloud and then added to the complete point cloud generated using the folding technique.

$$f^j(x) = \frac{\sum_{i=1}^k w_i(x) f_i^j}{\sum_{i=1}^k w_i(x)} \quad (6)$$

$$w_i(x) = \frac{1}{d(x, x_i)} \quad (7)$$

The interpolation operation is shown in (6), where $f^j(x)$ is the interpolated feature from the lower level (coarse point cloud or middle point cloud). f_i^j are features of the k -nearest neighbors of the higher level (middle point cloud or complete point cloud) in the lower level. Each neighbor is multiplied with a weight value which is equal to the distance of the point in a higher level to the lower level. After generating the middle point cloud using the folding technique, interpolate the coarse point cloud using the middle point cloud and add them, similarly after generating the complete point cloud, interpolate the middle point cloud using the complete point cloud and add them. The number of nearest neighbors used to interpolate is 3. The interpolation technique shown in (6) is proposed in PointNet++ [14]. In [14], it is used for segmentation of the input point cloud.

3.3 Capsule Point Completion Network

A. Capsule based Decoder:

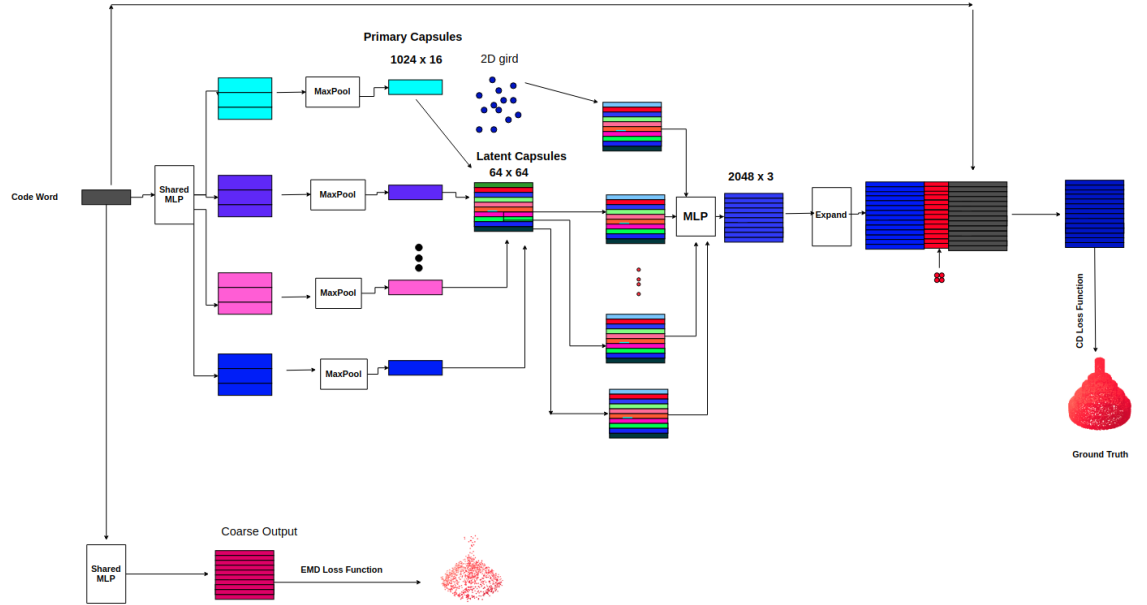


Figure 24 Capsule based decoder architecture.

Dynamic routing between capsules [47] has been shown to have effective performance on the MNIST dataset with very shallow architectures. [52] extended the dynamic routing architecture for multiple layers. [46] is the first architecture to apply the concept of dynamic routing to 3D point clouds. They have shown impressive results with reconstructing the input point cloud. Our capsule-based architecture as shown in Figure 21 uses the dynamic routing algorithm between a primary point capsule and latent capsules to obtain the final complete point cloud.

The primary point capsules are produced by applying a different shared multi-layer perceptron and then maxpooling operation to the codeword obtained from the encoder. To compute the latent capsules, the dynamic routing algorithm is applied to the primary point capsules. Applying the dynamic routing algorithm makes the model attend to different parts of the coarse output. The latent capsules are concatenated with the tiled fixed 2D grid and then applied shared multi-layer perceptron to compute the final complete point cloud. The model is optimized in stages of coarse output and final complete point cloud.

$$V_j = \frac{\|S_j\|^2}{1 + \|S_j\|^2} \frac{S_j}{\|S_j\|} \quad (7)$$

In (7) V_j is the vector output of capsule j and S_j is its total input.

$$S_j = \sum_i c_{ij} U_{ij} \quad (8)$$

In (8) j denotes the latent capsule and i denotes primary capsule. Input to a latent capsule is the prediction vectors of primary capsule multiplied by coupling coefficient c_{ij} .

$$U_{ij} = W_{ij} u_i \quad (9)$$

Equation (9) shows the prediction vector of a primary capsule is obtained by multiplying output of primary capsule u_i with the weight matrix.

$$c_{ij} = \frac{\exp(b_{ij})}{\sum_k \exp(b_{ik})} \quad (10)$$

Equation (10) shows how coupling coefficients are computed. The coupling coefficients from primary capsule i to latent capsule j sum to 1 and are computed by (10) which is called routing softmax, whose initial logits b_{ij} are the log prior probabilities that primary capsule i connected to upper latent capsule j .

B. Training Strategy

The model is coded in the TensorFlow framework and trained for 10 epochs with an initial learning rate of 0.001. The learning rate is decayed after every 1000 iterations with decay rate of 0.1. For the coarse output loss, the EMD [23] loss function is used and for the complete point cloud, the CD [23] loss function is used. The number of iterations for calculating coupling coefficients in each training iteration are three.

3.4 Multiview Point Completion Network

A. Multiview Encoder:

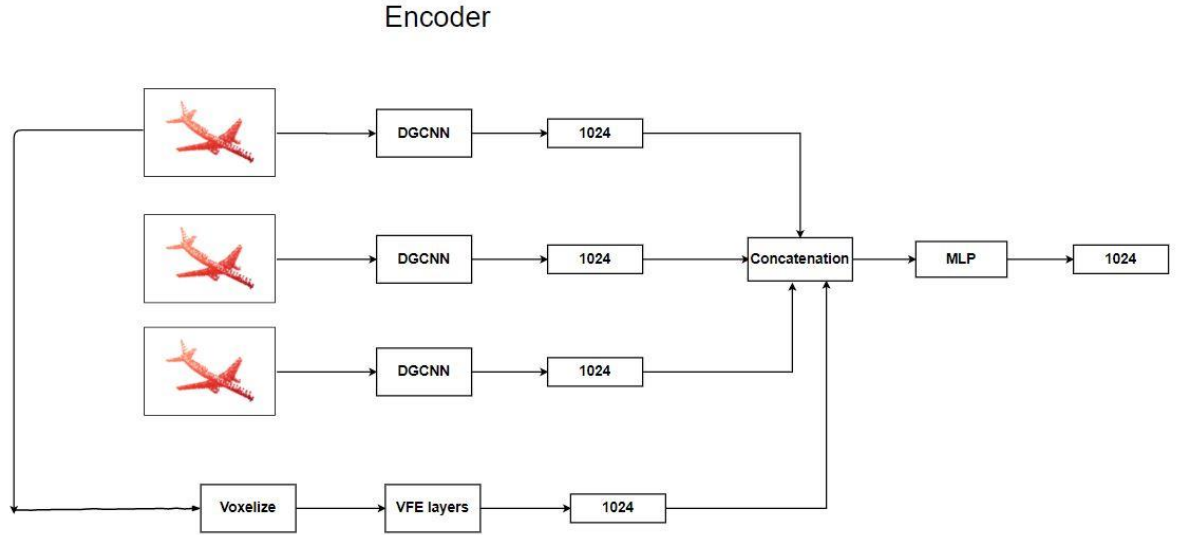


Figure 25 Multiview encoder.

The encoder of the multiview point completion network using DGCNN [16] is shown in Figure 22. The input partial point cloud is rotated by 1 degree and 3 degrees. All three point clouds, the input point cloud without rotation and the two rotated point clouds are fed into DGCNN [16] to compute the global vector of size 1024. To construct the k-nearest neighbors in DGCNN [16], we used $k=20$ and computed the edge features as in [16]. The input partial point cloud is voxelized and computed voxel-wise features using voxel feature extraction (VFE) layers, and applied maxpooling to compute the global vector of size 1024. We used three stacked VFE layers to extract the voxel-wise features. In general, VFE layer computes point-wise features and then applies symmetric maxpooling function to obtain the global voxel-wise features. In this way, we applied 3 VFE layers after converting the input partial point cloud to voxels. These global voxel-wise features are concatenated with the point-wise features to obtain the final voxel-wise features. All the global vectors computed from the rotated point clouds and VFE layers are fed into a multi-layer perceptron to compute the final codeword of size 1024. The multi-layer perceptron architecture has $512 - 1024$ nodes to compute the final codeword.

B. Decoder

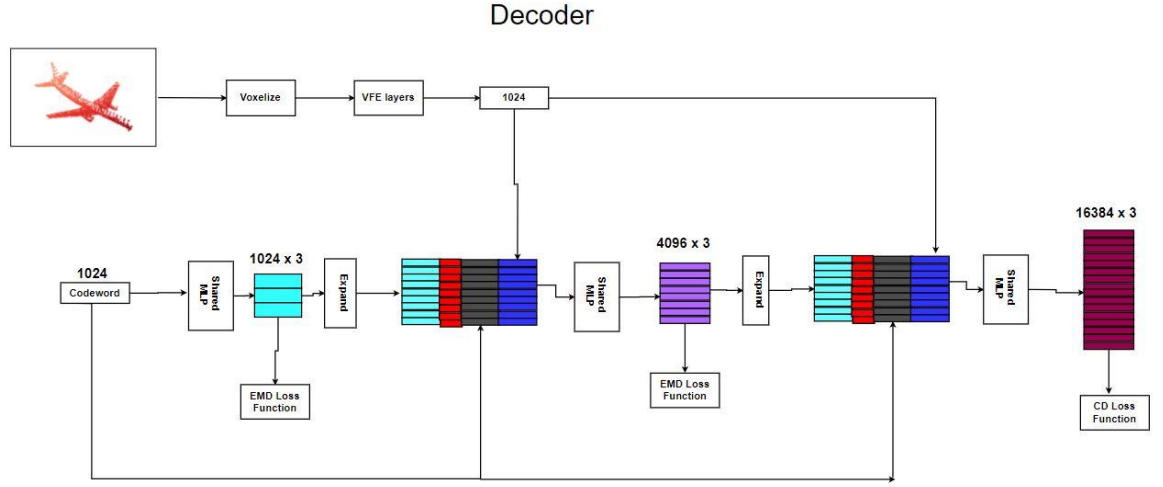


Figure 26 Decoder for Multiview point completion network.

The decoder for the multiview point completion network is shown in Figure 23. The decoder receives a codeword and a global voxel-wise feature vector from the encoder and optimizes it in multiple stages. In the first stage, the codeword is converted into a coarse point cloud of shape 1024×3 by applying a shared multi-layer perceptron. The loss is calculated using the EMD [23] loss function. In the second stage, a 2D random grid is concatenated to the coarse vector obtained in the first stage. The codeword is tiled and concatenated and global voxel-wise feature is concatenated to obtain a middle vector of shape 4096×2053 . The shared multi-layer perceptron is applied to this vector and the loss is calculated by comparing with the ground truth using EMD [23] loss function. In the third stage, the process is repeated as in the second stage, whereby a 2D random grid is concatenated to the middle vector, a coarse vector obtained in the first stage is tiled and concatenated, and the codeword from the encoder is tiled and concatenated with a voxel-wise global feature vector. Then shared multi-layer perceptron is applied to compute the final complete point cloud.

The ground truth for the coarse point cloud and the middle point cloud is generated by using the farthest sampling algorithm [14] from the complete point cloud ground truth. Iterative farthest sampling is used since it is better than random sampling. The coarse point cloud and the middle point cloud is optimized using the EMD [23] loss function, as the cost of EMD [23] is higher than of the CD [23] loss function.

C. Training Strategy

The model is coded in the TensorFlow framework and trained for 10 epochs with an initial learning rate of 0.001. The learning rate is decayed after every 1000 iterations with a decay rate of 0.1. For the coarse output loss, the EMD [23] loss function is used, and for the complete point cloud, the CD [23] loss function is used.

3.5 Self-supervised based point completion network

A. Self-supervised Encoder:

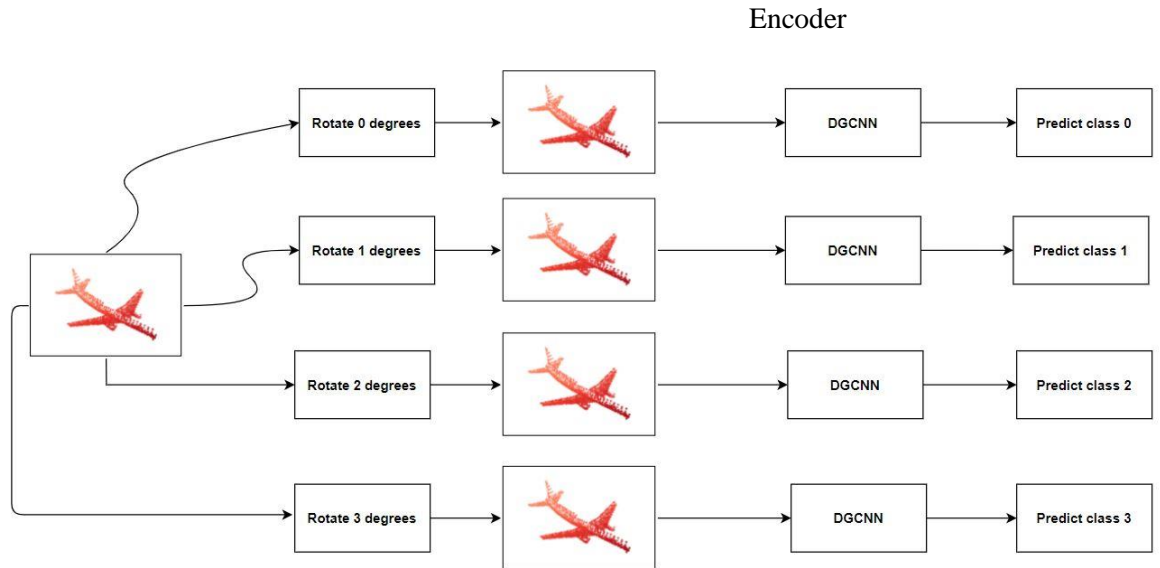


Figure 27 Self-supervised pretext task for point completion network.

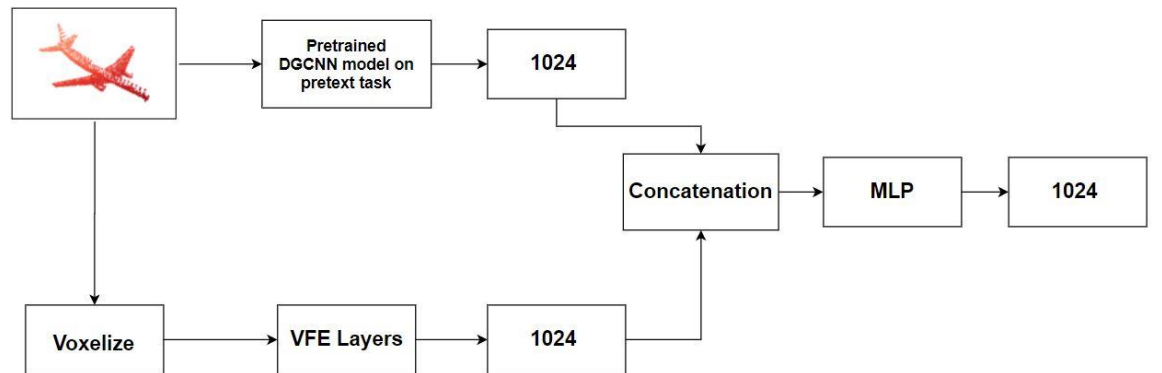


Figure 28 Encoder for self-supervised learning.

The encoder for the self-supervised point completion network is shown in Figure 25. We rotated the input partial point cloud by angles of 0 , 1 , 2 and 3 degrees, then trained the DGCNN [16] network to classify the rotation class it belongs. This is the pretext task which we have chosen to train the model in a self-supervised way. After training the model for this rotation pretext task, we perform transfer learning using the pretext weights to compute the global vector of size 1024 . This global vector is concatenated with the global voxel-wise feature computed using voxel feature extraction layers and then a shared multi-layer perceptron is applied to compute the final codeword of size 1024 . The model trained on pretext task learns the semantic features of the input partial point clouds in an improved fashion. Using transfer learning with those weights extracts the features of the input partial clouds with better results. To classify the rotated point clouds, the model must identify certain features which helps in correctly classifying the point clouds rotated in different directions. This identification of features by the model for this pretext task helps the task of point completion with less data.

B. Decoder:

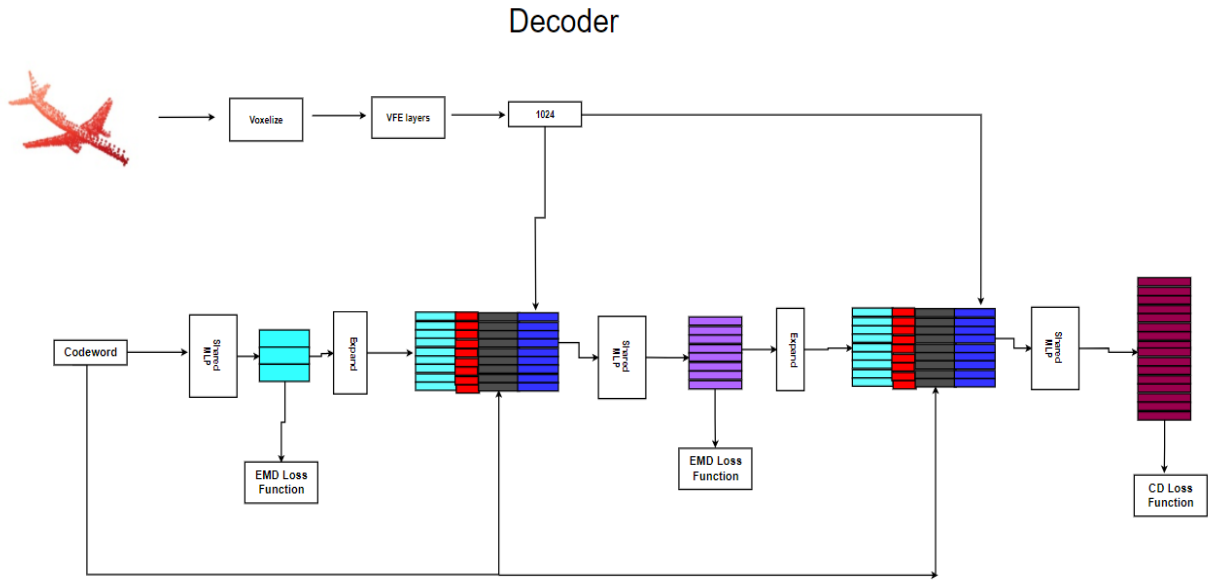


Figure 29 Decoder for self-supervised learning.

The decoder for self-supervised point completion network is shown in Figure 26. A coarse vector of shape 1024×3 is computed using the codeword from the encoder, and then applying a shared multi-layer perceptron. The coarse vector is optimized using the EMD loss function. A 2D grid is concatenated to the coarse vector along with the tiled codeword and global voxel-wise feature

vector, then a shared-multi layer perceptron is applied to compute the middle vector of shape 4096×3 . The middle vector is optimized using the EMD loss function. The middle vector is concatenated with a 2D grid, a tiled code word, and a tiled voxel-wise global feature vector to compute the fine vector of shape 16384×3 . This fine vector is the complete point cloud prediction of the model. The fine vector is optimized using the CD loss function. Since the coarse vector and middle vector are smaller shapes, we used the EMD loss function for those and the CD loss function which is less compute heavy compared to EMD is used for the final fine vector.

C. Training Strategy

The model is coded in the TensorFlow framework and trained for 10 epochs with an initial learning rate of 0.001. The learning rate is decayed after every 1000 iterations with decay rate of 0.1. For the coarse output loss, the EMD [23] loss function is used and for the complete point cloud, the CD [23] loss function is used.

3.6 Mesh Segmentation

For the robot gripper like a suction cup to hold the object firmly, it has to identify the good faces on the surface of the object. Since 3D data can be represented in different ways, representing the 3D object in the form of 3D mesh provides more information about the connectivity of faces on the surface of the object. We represented 3D objects in the form of a 3D mesh and apply deep learning based architectures to identify the good faces on the object. This problem is formulated as mesh segmentation, where each face in the 3D mesh can be a good or bad face.

A. PointNet Center:

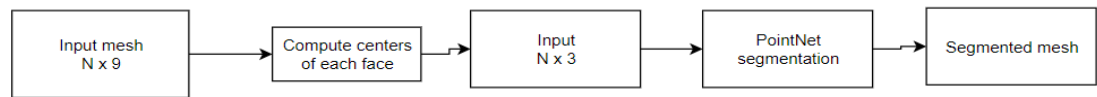


Figure 30 PointNet Center.

The PointNet center [15] architecture is shown in Figure 27. Input to the architecture is a 3D mesh. The mesh consists of faces and vertices. The faces contain information of connection between vertices and vertices are actual 3D points of the mesh. We computed the centers of each face in the 3D input mesh of shape $N \times 9$ to convert into shape of $N \times 3$. Then PointNet [15] segmentation architecture is applied to this modified input to output segmented mesh.

B. PointNet Mesh:

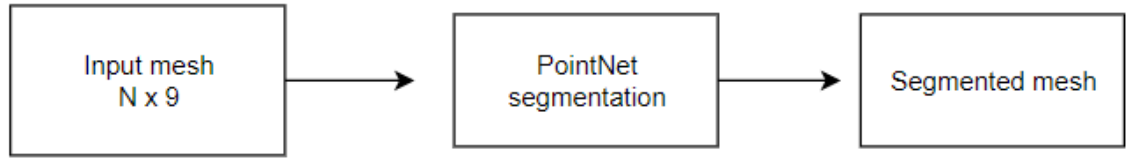


Figure 31 PointNet Mesh.

The PointNet mesh architecture is shown in Figure 28. PointNet [15] uses the 3D point clouds of shape $N \times 3$ as inputs, whereas the architecture shown in Figure 28, uses 3D mesh as input. We concatenated the vertices of each face in the 3D mesh to convert into a vector of shape $N \times 9$. Then segmentation architecture in [15] is applied to segment the input mesh.

C. DGCNN Center:

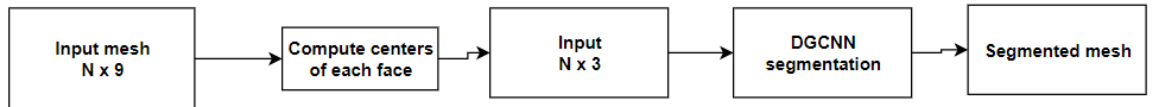


Figure 32 DGCNN Center.

DGCNN [16] architecture we used for mesh segmentation is shown in Figure 29. As shown in Figure 29, we computed the centers of each face in the input 3D mesh of shape $N \times 9$ and converted into input of shape $N \times 3$. The DGCNN [16] architecture is applied to segment the 3D input mesh. The number of segmentation classes in the output is two, one is good face and the other is bad face. We used the cross entropy loss in the output layer to calculate the loss between the predicted and ground truth.

D. MeshNet

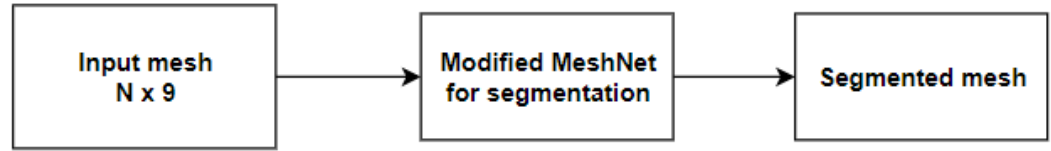


Figure 33 MeshNet for segmentation.

MeshNet [48] considers the face information, normal vectors of each face, neighboring corners of each face and centers of each face to perform the classification task. We processed the input 3D mesh in the same way as in [48], computed neighboring corners of each face, normal vectors of each face and centers of each face. We modified [48] for segmentation and used same architecture for the input 3D mesh

Chapter 4

Implementation

4.1 Datasets

A. ShapeNet

ShapeNet[8] is a large-scale annotated repository of 3D CAD models of different objects. It provides a rich set of annotations for every object. ShapeNet provides variety of 3D models to evaluate the performance of 3D shape reconstruction, 3D part segmentation and 3D object classification. There are total 55 common categories of 3D objects in ShapeNet.

ID	Name	Num	ID	Name	Num	ID	Name	Num
04379243	table	8443	03593526	jar	597	04225987	skateboard	152
02958343	car	7497	02876657	bottle	498	04460130	tower	133
03001627	chair	6778	02871439	bookshelf	466	02942699	camera	113
02691156	airplane	4045	03642806	laptop	460	02801938	basket	113
04256520	sofa	3173	03624134	knife	424	02946921	can	108
04090263	rifle	2373	04468005	train	389	03938244	pillow	96
03636649	lamp	2318	02747177	trash bin	343	03710193	mailbox	94
04530566	watercraft	1939	03790512	motorbike	337	03207941	dishwasher	93
02828884	bench	1816	03948459	pistol	307	04099429	rocket	85
03691459	loudspeaker	1618	03337140	file cabinet	298	02773838	bag	83
02933112	cabinet	1572	02818832	bed	254	02843684	birdhouse	73
03211117	display	1095	03928116	piano	239	03261776	earphone	73
04401088	telephone	1052	04330267	stove	218	03759954	microphone	67
02924116	bus	939	03797390	mug	214	04074963	remote	67
02808440	bathtub	857	02880940	bowl	186	03085013	keyboard	65
03467517	guitar	797	04554684	washer	169	02834778	bicycle	59
03325088	faucet	744	04004475	printer	166	02954340	cap	56
03046257	clock	655	03513137	helmet	162			
03991062	flowerpot	602	03761084	microwaves	152		Total	57386

Figure 34 Different categories for ShapeNet dataset [8].

Figure 31 shows the different categories and their corresponding dataset size in ShapeNet.



Figure 35 Examples of ShapeNet dataset [8].

Figure 32 shows some examples of 3D CAD objects from ShapeNet.

B. Point Completion Networks:

For the point completion task, we need a dataset containing pairs of partial and complete points clouds. We used the same method as in PCN [2] to generate complete and partial point clouds. The complete point clouds are generated by sampling 16,384 points from mesh surfaces of 3D CAD models of the ShapeNet dataset. To generate partial point clouds, we back projected 2.5D depth images into 3D. There are total of eight categories of objects chosen for this task: airplane, cabinet, chair, car, lamp, sofa, table, vessel. Back-projecting the depth images into 3D brings the data distribution to the real-world sensor data. Since real world data won't contain detailed information of 3D objects, synthetic 3D objects like those in ShapeNet provide rich 3D information to process.



Figure 36 Examples of PCN dataset [2].

Examples of dataset containing input partial point clouds and ground truth complete point clouds are shown in Figure 33.

C. Phenix Automation:

For the robot grippers like suction cups to hold the objects firmly, the 3D mesh object should contain a minimum flat surface area. Since a gripper can hold the object anywhere, a segmentation of the 3D mesh satisfying the minimum area requirements should be done. Present datasets like ShapeNet [], ScanNet [49], S3DIS [50] do not contain 3D objects which would typically might be found in industrial machinery environments. As such, we have created our own dataset grabbing different 3D CAD models from GrabCAD [51] and then manually annotated the dataset. For annotating the dataset, we developed a MATLAB tool.

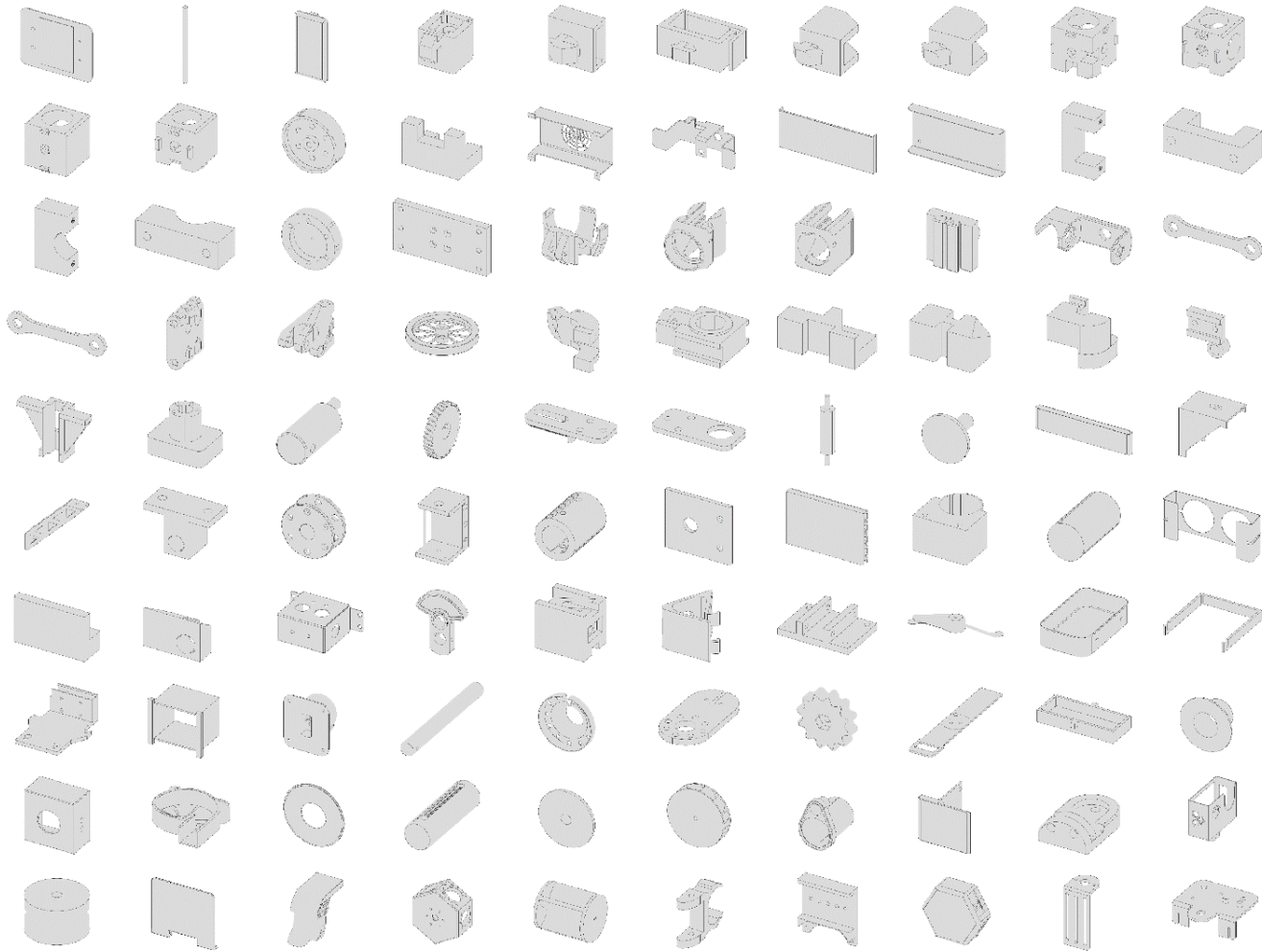


Figure 37 Examples of the Phenix Automation dataset.

3D data can be represented in many forms like point clouds, voxels, and mesh. But for the task of the robot gripper, representing the 3D object in the form of point clouds will not help the task in hand. Since point clouds are sparse, it just provides the information about various 3D points and little about the relationship between those points. To use point clouds for this task, the relationship between different points should be discovered and then points suitable for robot gripping should be identified. There exists ambiguity in the case of edges. For example, should a point on the edge be considered for robot gripping? Similar problems arise if 3D data is represented in the form of voxels, since each voxel just stores the information of 3D points belonging to that voxel, and as such, this information may not be useful for this task. Representing 3D data in the form of a 3D surface mesh provides more information of the 3D object like faces and neighboring

connectivity. Since faces provide rich surface information of the 3D object, using the information of faces can help identify the good faces for our robot gripper's suction cup to adhere to.

Examples of different 3D CAD models fetched online are shown in Figure 34. CAD models fetched meet the real-world distribution of different industrial machinery parts. We generated 3D mesh from the CAD model using the MATLAB PDE toolbox, and annotated good and bad faces. We created a MATLAB tool for faster dataset annotation.

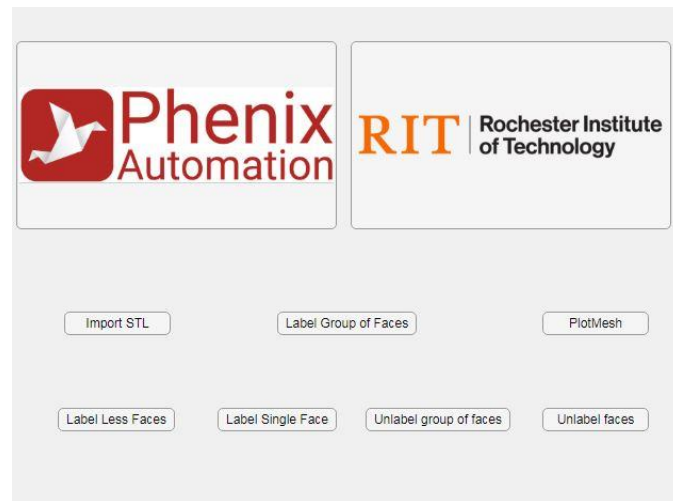


Figure 38 View of MATLAB tool for dataset annotation.

The MATLAB tool developed for dataset annotation purpose is shown in Figure 35. All the 3D CAD models fetched online are converted into the MATLAB supported format STL (Stereolithography). We used the PDE (Partial Differential Equations) toolbox in MATLAB to generate tetrahedral 3D mesh surfaces, since we are interested only on outside body surfaces. We converted tetrahedral 3D mesh into triangular 3D meshes and annotated good vs. bad faces using MATLAB callback functions.

In the MATLAB tool shown in Figure 35, after we import the STL file, the tool will automatically generate a triangular 3D mesh of the 3D object, and then user can annotate all the good faces of the 3D object. A user can label group of faces or single depending on the requirement.

Examples of training data are shown in Figure 36. The color red corresponds to bad faces and the color cyan corresponds to good faces. A good face should be a minimum area of 10 mm both on X and Y axis. Some examples in Figure 36 don't contain any good faces for the robot gripper to hold. All the training data has different shapes. Some parts contain more triangular elements (faces) and some parts contain lesser triangular elements.

Training Data

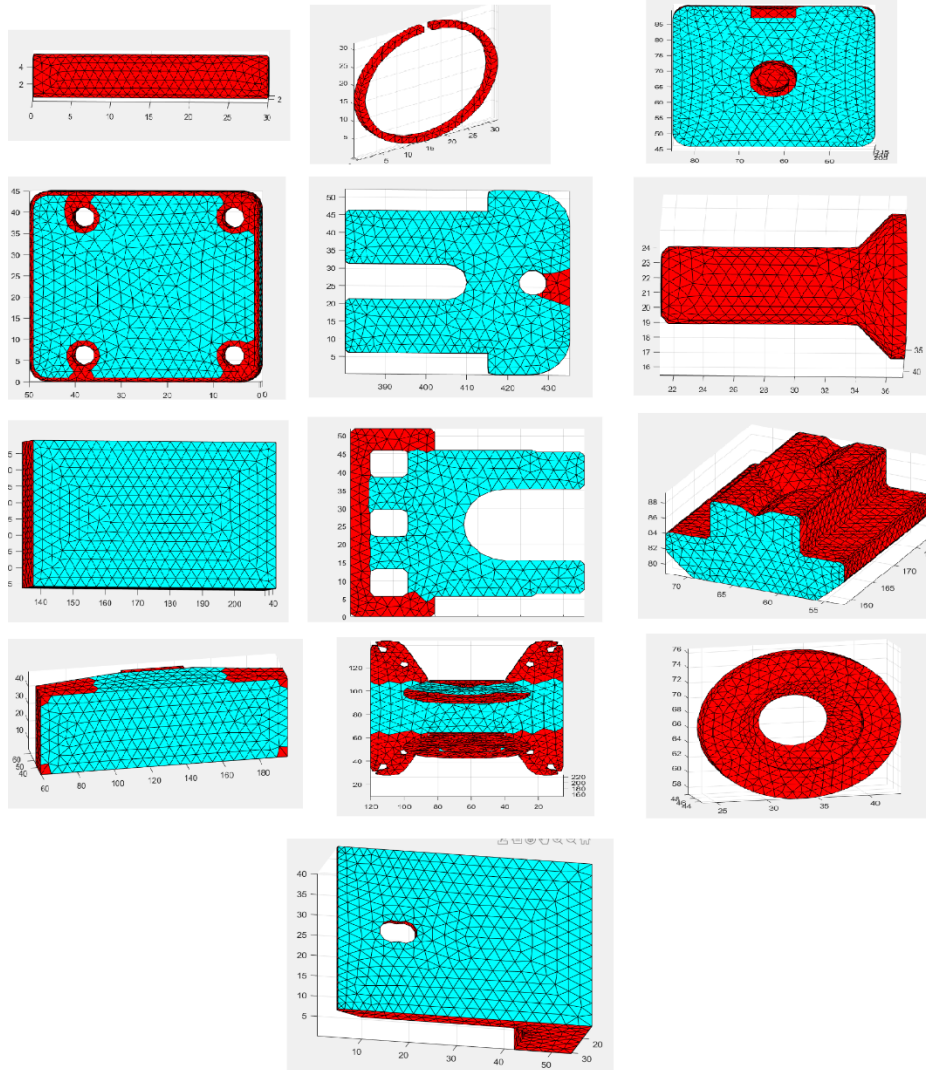


Figure 39 Examples of annotated dataset for robot gripper.

Chapter 5

Results and Analysis

5.1 Results

A. MS PCN:

Table 1 Results of MS PCN.

Model	CD (Chamfer Distance)
PCN (Baseline)	0.00986
MS PCN	0.013018
MS PCN (stagewise)	0.01015
MS PCN (no voxels)	0.0150

In Table 1, the results are shown for the multistage point completion network (MS PCN). Our results are compared with the baseline results. Metric compared is the average chamfer distance loss of the test dataset. MS PCN average loss is high compared to the baseline results. It seems that voxel features and multistage optimization of the decoder isn't helping the model to output a complete point cloud. Table 1 shows the results of stagewise training as well, results of stagewise training and normal training are almost the same. But loss can be improved definitely with the stagewise training, since in optimizing the middle point cloud from coarse point cloud the network is smaller and data is huge, model is not able to produce optimal middle point cloud equivalent to the ground truth point cloud. Similarly, the trained network may be smaller to optimize the complete point cloud from the middle point cloud.

Table 1 shows the results of the MS PCN model without voxels as well, in this experiment we didn't voxelized the input partial point cloud, we used the similar encoder architecture as in PCN [2]. But loss is high compared to the models with voxel features, so adding voxel features definitely helping the model to learn some features for completing the point cloud. We trained our models with different hyperparameters like changing the learning rate schedule, initial learning rate, different optimizers, and different batch sizes, we provided our best results obtained during training.

B. Edge convolution-based point completion network:

Table 2 Results of Edge convolution-based PCN.

Model	CD (Chamfer Distance)
PCN (Baseline)	0.00986
Edge convolution PCN	0.009368

Table 2 shows the results of the edge convolution-based encoder and interpolated decoder. Since edge convolution extracts the neighborhood information much better than the baseline model because in the edge convolution a k nearest neighbor graph is constructed this helps the model to effectively extract the local neighborhood compared to other models. And the final output of the voxel convolution also concatenated to the edge convolution output. As voxel convolution does the interaction between points in different voxels, it helps in understanding the overall shape of the input partial point cloud. In the interpolated decoder instead of replicating directly coarse point cloud and middle point cloud, an interpolated coarse point cloud and an interpolated middle point cloud are concatenated. As directly replicating the coarse point cloud and the middle point cloud doesn't much information to the further stages, an interpolated version adds more information since we are optimizing the coarse point and the middle point cloud using earth mover distance (EMD) loss function. In this optimization stage, an interpolated coarse point cloud and the middle point cloud can optimize much better compared to the non-interpolated decoders. And the final complete point cloud is optimized using the CD loss function. As an effect of all these architectural changes, we can see an improvement in the loss compared to the baseline model.

C. Capsule based decoder:

Table 3 Results of Capsule based decoder.

Model	CD (Chamfer Distance)
PCN (Baseline)	0.00986
Capsule based decoder	0.02061

Table 3 shows the results of the point completion network with the capsule-based decoder. The loss of the model with the capsule-based decoder is high compared to the baseline results. In the capsule-based model, we are optimizing the coarse point cloud by applying a multilayer perceptron on the codeword from the encoder and the complete point cloud by employing the dynamic routing algorithm. Because of this multitasking setup, maybe the weights learned using the dynamic routing algorithm are getting disturbed. And the final complete point cloud is not generated by an end-to-

end dynamic routing algorithm from the codeword, initially, an intermediate point cloud of shape 2048×3 is obtained by dynamic routing algorithm and the complete point cloud is obtained by folding operation from the intermediate point cloud.

D. Multiview PCN:

Table 4 Results of Multiview PCN.

Model	CD (Chamfer Distance)
PCN (Baseline)	0.00986
DGCNN Multiview	0.0176

The results of the multiview PCN are shown in Table 4. The input to the encoder of the multiview PCN is the input partial clouds rotated by 0 degrees, 1 degree and 3 degree, and used edge convolution to extract the features, concatenated all the features with the voxel features. The decoder is same as the MS PCN model, our results unable to outperform the baseline results. Rotating the point cloud by different angles didn't help the model to learn new features to complete the point cloud, model is able to learn the same features in the rotated point clouds even with edge convolutions.

a. Self-supervised model:

Table 5 Results of Self-supervised PCN.

Model	CD (Chamfer Distance)
PCN (Baseline)	0.00986
Self-supervised model	

E. PointNet-Center:

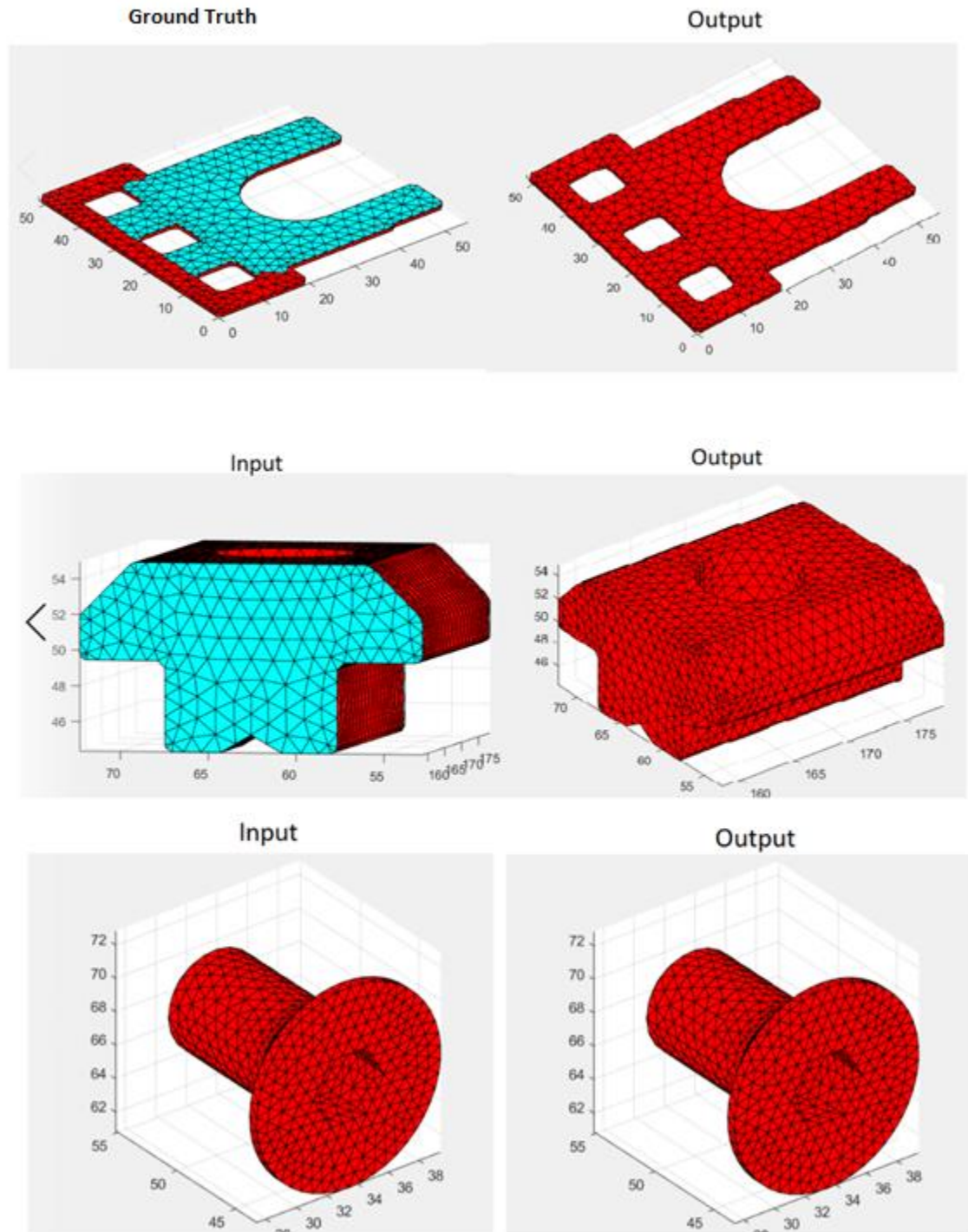


Figure 40 Results of PointNet-Center.

The results of PointNet-Center for 3D mesh segmentation is shown in the Figure 40. The dataset used to train the model is shown in Figure 39. The input to the model is the centers of each mesh in the input dataset, model is not able to distinguish between good face and bad face, it always outputs every face of the mesh as bad for gripping. In the Figure 40, the ground truth and the corresponding outputs are shown. Many for the input samples having good faces, the model outputs every face as bad face. Since the PointNet [15] architecture just learns the pointwise features, this information alone is not sufficient to segment the input 3D mesh.

F. PointNet-Mesh:

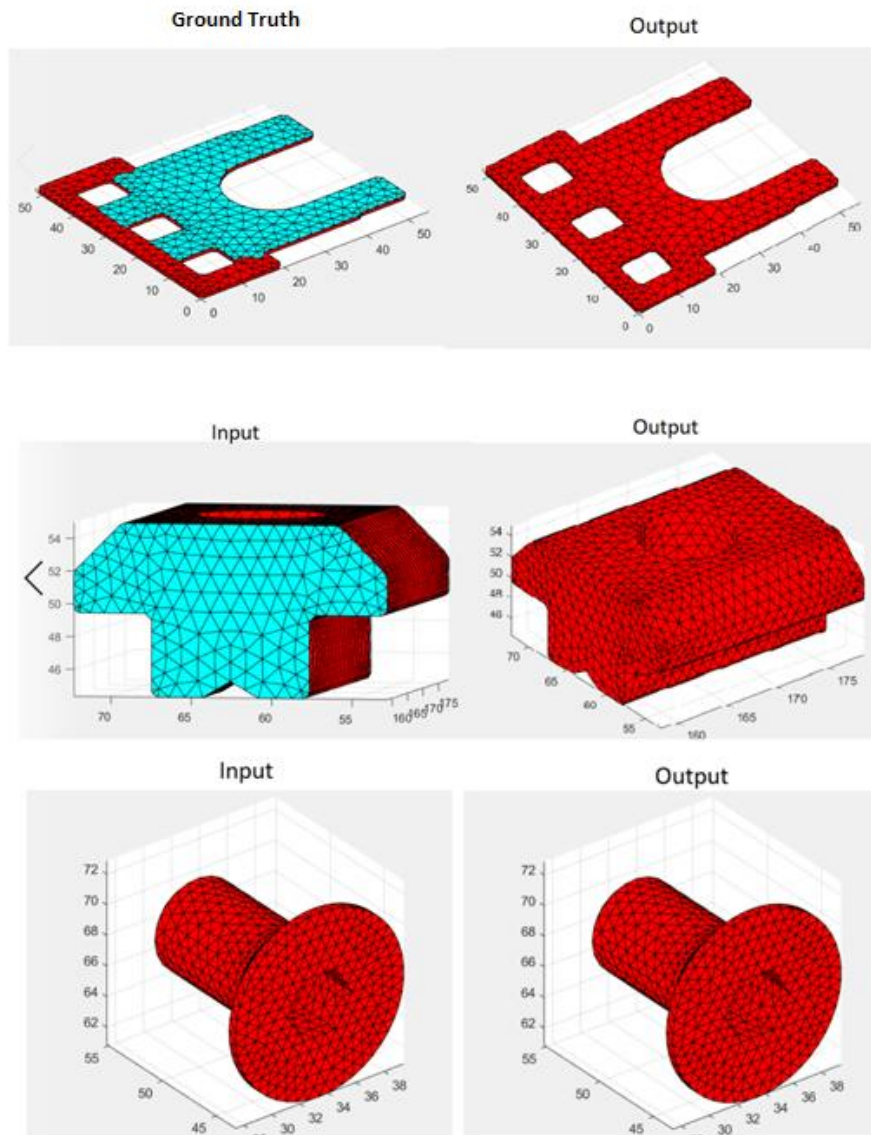


Figure 41 Results of PointNet-Mesh.

The results of 3D mesh segmentation for the model PointNet-Mesh is shown in Figure 41. The input to the model is concatenated vector of all corners in each face of the 3D input mesh. Even though the input has many samples having good faces, the model always outputs each face as a bad face. Even though the input has information about each mesh, but the PointNet [15] style architecture unable to learn the relation between the corners of each face and different faces in the input 3D mesh.

G. DGCNN-Center:

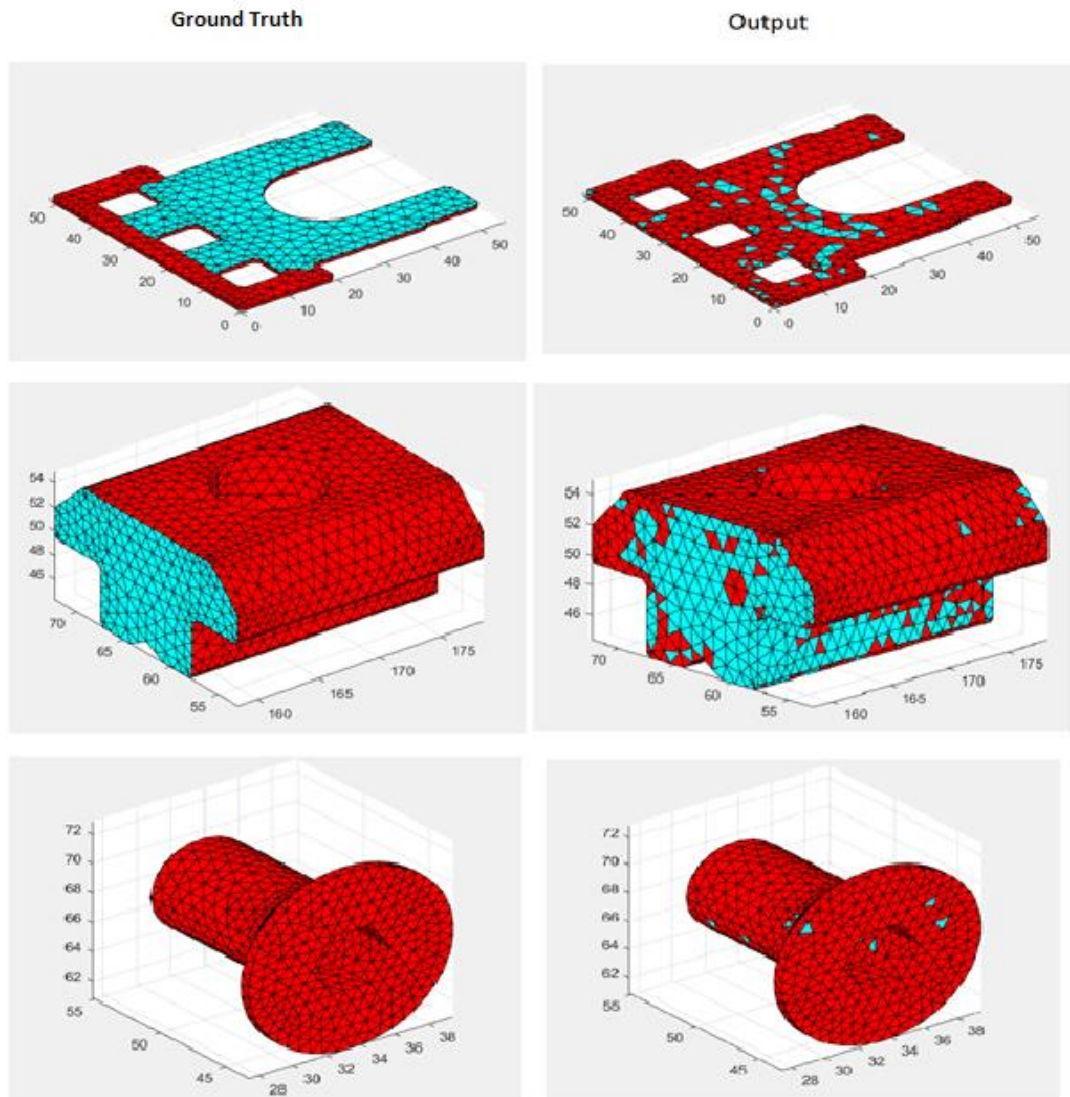


Figure 42 Results of DGCNN-Center.

The results of the model DGCNN-Center are shown in Figure 42. The input to the model is the centers of each mesh obtained from the 3D mesh shown in Figure 39. Unlike PointNet [15] style models, DGCNN [16] based model able to output true good faces. Even though the dataset is small, the model is able to learn the relation between different faces in the 3D mesh to classify good and bad faces. Since [16] constructs a k -nearest graph, it extracts the local neighborhood information much better than PointNet [15] style based architectures.

H. MeshNet:

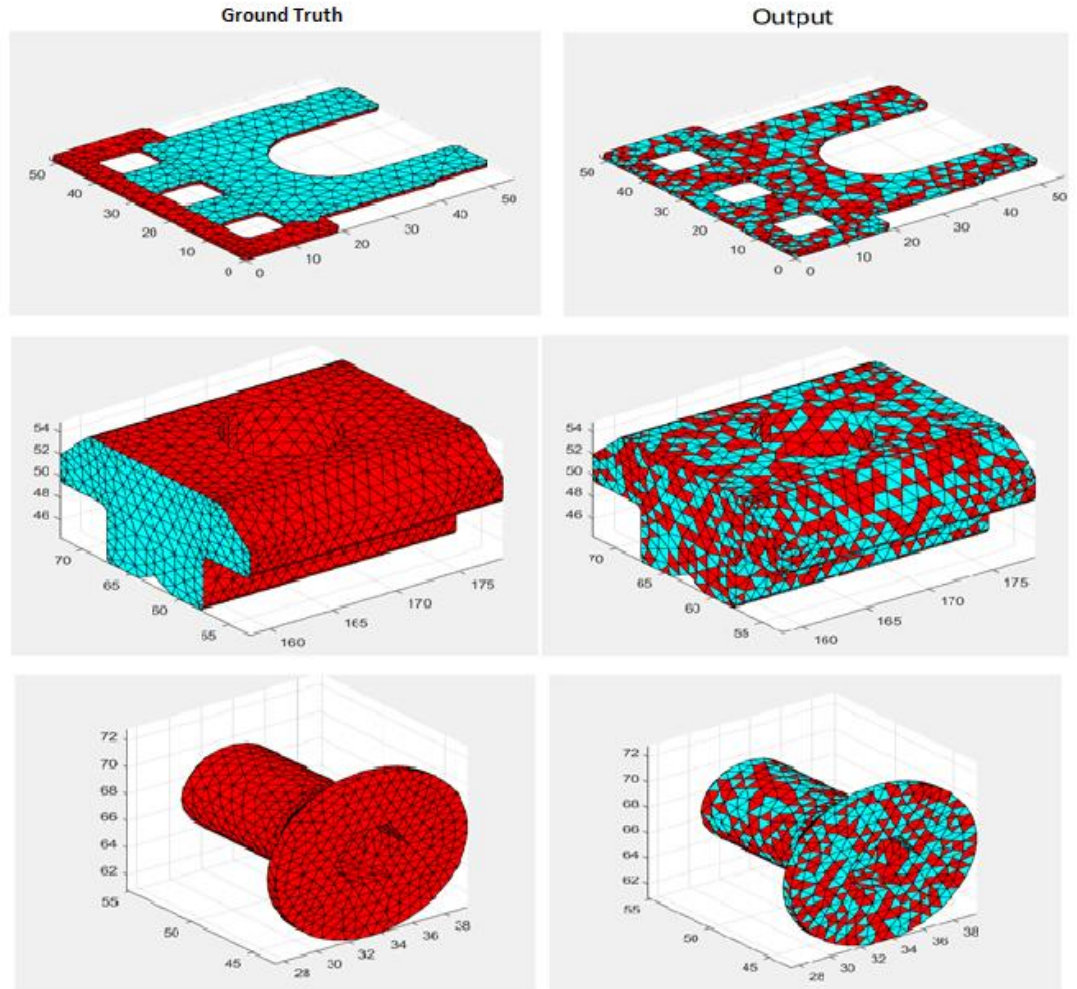


Figure 43 Results of MeshNet.

The results of the MeshNet is shown in Figure 43. We used the same architecture as in [48], the input to the model contains face information, normal of each face and neighbors of each face. With the limited training dataset, the model is more inclined to output most of the faces as a good face. In Figure 43, we can see even for the input sample which doesn't contain any good face, the model predicted many good faces. DGCNN-Center has better results compared to the MeshNet, even though MeshNet considers each face as the operating unit and applies mesh convolutions, DGCNN-Center model able to learn the relation between different faces with just information about the centers.

Bibliography

- [1] D.-Y. Chen, X.-P. Tian, Y.-T. Shen, and M. Ouhyoung. On visual similarity based 3d model retrieval. In *Computer graphics forum*, volume 22, pages 223–232. Wiley Online Library, 2003.
- [2] Yuan, Wentao, et al. "Pcn: Point completion network." 2018 International Conference on 3D Vision (3DV). IEEE, 2018.
- [3] H. Su, S. Maji, E. Kalogerakis, and E. G. Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition.
- [4] R. B. Rusu, N. Blodow, and M. Beetz. Fast point feature histograms (fpfh) for 3d registration. In *Robotics and Automation, 2009. ICRA'09. IEEE International Conference on*, pages 3212–3217. IEEE, 2009.
- [5] Y. Li, R. Bu, M. Sun, and B. Chen. Pointcnn. arXiv preprint arXiv:1801.07791, 2018.
- [6] M. Aubry, U. Schlickewei, and D. Cremers. The wave kernel signature: A quantum mechanical approach to shape analysis. In *Computer Vision Workshops (ICCV Workshops)*, 2011 IEEE International Conference on, pages 1626–1633. IEEE, 2011.
- [7] X. Han, Z. Li, H. Huang, E. Kalogerakis, and Y. Yu. Highresolution shape completion using deep neural networks for global structure and local geometry inference. arXiv preprint arXiv:1709.07599, 2017.
- [8] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1912–1920, 2015.
- [9] D. Stutz and A. Geiger. Learning 3d shape completion from laser scan data with weak supervision.
- [10] C. R. Qi, H. Su, M. Nießner, A. Dai, M. Yan, and L. Guibas. Volumetric and multi-view cnns for object classification on 3d data. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2016.
- [11] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected networks on graphs. arXiv preprint arXiv:1312.6203, 2013.
- [12] J. Masci, D. Boscaini, M. Bronstein, and P. Vandergheynst. Geodesic convolutional neural networks on riemannian manifolds. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 37–45, 2015.
- [13] H. Ling and D. W. Jacobs. Shape classification using the inner-distance. *IEEE transactions on pattern analysis and machine intelligence*, 29(2):286–299, 2007.

- [14] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In Proc. NIPS, 2017.
- [15] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In Proc. CVPR, 2017.
- [16] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon. Dynamic graph cnn for learning on point clouds. arXiv preprint arXiv:1801.07829, 2018.
- [17] R. B. Rusu, N. Blodow, Z. C. Marton, and M. Beetz. Aligning point cloud views using persistent feature histograms. In 2008 IEEE/RSJ International Conference on Intelligent Robots and Systems, pages 3384–3391. IEEE, 2008.
- [18] H. Su, V. Jampani, D. Sun, S. Maji, E. Kalogerakis, M.-H. Yang, and J. Kautz. Splatnet: Sparse lattice networks for point cloud processing. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 2530–2539, 2018.
- [19] D. Maturana and S. Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In IEEE/RSJ International Conference on Intelligent Robots and Systems, September 2015.
- [20] J. Varley, C. DeChant, A. Richardson, J. Ruales, and P. Allen. Shape completion enabled robotic grasping. In Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on, pages 2442–2447. IEEE, 2017.
- [21] B. Yang, S. Rosa, A. Markham, N. Trigoni, and H. Wen. 3d object dense reconstruction from a single depth view. arXiv preprint arXiv:1802.00411, 2018.
- [22] P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. Guibas. Learning representations and generative models for 3d point clouds.
- [23] H. Fan, H. Su, and L. Guibas. A point set generation network for 3d object reconstruction from a single image. In Conference on Computer Vision and Pattern Recognition (CVPR), volume 38, 2017.
- [24] J. Sun, M. Ovsjanikov, and L. Guibas. A concise and provably informative multi-scale signature based on heat diffusion. In Computer graphics forum, volume 28, pages 1383–1392. Wiley Online Library, 2009.
- [25] A. E. Johnson and M. Hebert. Using spin images for efficient object recognition in cluttered 3d scenes. IEEE Transactions on pattern analysis and machine intelligence, 21(5):433–449, 1999.
- [26] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet Classification with Deep Convolutional Neural Networks,” in Advances in Neural Information Processing Systems 25, F.

Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 1097–1105.

[27] C. Szegedy et al., “Going Deeper With Convolutions,” presented at the Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 1–9.

[28] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” presented at the Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.

[29] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, “Improving neural networks by preventing co-adaptation of feature detectors,” arXiv:1207.0580 [cs], Jul. 2012.

[30] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In CVPR09, 2009.

[31] Karen Simonyan and Andrew Zisserman, “Very deep convolutional networks for large-scale image recognition,” arXiv preprint arXiv:1409.1556, 2014.

[32] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation,” presented at the Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2014, pp. 580–587.

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

[34] Yang, Yaoqing, et al. "Foldingnet: Point cloud auto-encoder via deep grid deformation." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2018.

[35] Gidaris, Spyros, Praveer Singh, and Nikos Komodakis. "Unsupervised representation learning by predicting image rotations." arXiv preprint arXiv:1803.07728 (2018).

[36] Zhou, Yin, and Oncel Tuzel. "Voxelnet: End-to-end learning for point cloud based 3d object detection." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2018.

[37] Yifan, Wang, et al. "Patch-based Progressive 3D Point Set Upsampling." arXiv preprint arXiv:1811.11286 (2018).

[38] Mandikal, P., and Babu, R. V. 2019. Dense 3d point cloud reconstruction using a deep pyramid network. In Winter Conference on Applications of Computer Vision (WACV).

[39] <https://peltarion.com/knowledge-center/documentation/modeling-view/build-an-ai-model/blocks/1d-convolution-block>.

[40] <https://github.com/PetarV-/TikZ/tree/master/2D%20Convolution>

- [41] <https://computersciencewiki.org/index.php/Max-pooling> / Pooling
- [42] https://software.intel.com/sites/products/documentation/doclib/daal/daal-user-and-reference-guides/daal_prog_guide/GUID-9B434D4F-C723-4191-9A88-69148C75A3F1.htm
- [43] <https://towardsdatascience.com/complete-guide-of-activation-functions-34076e95d044>
- [44] <https://www.oreilly.com/library/view/tensorflow-for-deep/9781491980446/ch04.html>
- [45] <https://jhui.github.io/2017/11/03/Dynamic-Routing-Between-Capsules/>
- [46] Zhao, Yongheng, et al. "3D point capsule networks." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2019.
- [47] Sabour, Sara, Nicholas Frosst, and Geoffrey E. Hinton. "Dynamic routing between capsules." Advances in neural information processing systems. 2017.
- [48] Feng, Yutong, et al. "MeshNet: mesh neural network for 3D shape representation." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 33. 2019.
- [49] Dai, Angela, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. "Scannet: Richly-annotated 3d reconstructions of indoor scenes." In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5828-5839. 2017.
- [50] Armeni, Iro, et al. "Joint 2d-3d-semantic data for indoor scene understanding." arXiv preprint arXiv:1702.01105 (2017).
- [51] <https://grabcad.com/>
- [52] LaLonde, Rodney, and Ulas Bagci. "Capsules for object segmentation." arXiv preprint arXiv:1804.04241 (2018).
- [53] Ioffe, Sergey, and Christian Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift." arXiv preprint arXiv:1502.03167 (2015).
- [54] Groueix, Thibault, et al. "A papier-mâché approach to learning 3d surface generation." Proceedings of the IEEE conference on computer vision and pattern recognition. 2018.
- [55] Lang, Alex H., et al. "Pointpillars: Fast encoders for object detection from point clouds." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2019.
- [56] Huang, Gao, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. "Densely Connected Convolutional Networks." In CVPR, vol. 1, no. 2, p. 3. 2017.
- [57] Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition." arXiv preprint arXiv:1409.1556 (2014).
- [58] Zhang, Richard, Phillip Isola, and Alexei A. Efros. "Colorful image colorization." European conference on computer vision. Springer, Cham, 2016.
- [59] Noroozi, Mehdi, and Paolo Favaro. "Unsupervised learning of visual representations by solving jigsaw puzzles." European Conference on Computer Vision. Springer, Cham, 2016.

[60] Pathak, Deepak, et al. "Context encoders: Feature learning by inpainting." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.