



## Object synthesis by learning part geometry with multi-tasks

Sangpil Kim<sup>a,\*</sup>, Hyung-gun Chi<sup>a</sup>, Karthik Ramani<sup>a</sup>

<sup>a</sup>Purdue University, West Lafayette IN, 47906, USA

### ARTICLE INFO

#### Article history:

Received August 10, 2019

**Keywords:** Deep learning, Multi-task learning, Object synthesis, Object classification, Object reconstruction, Volumetric representation, Surface representation, Geometric

### ABSTRACT

We propose a conditional generative model, named Part Geometry Network (PG-Net), which synthesizes realistic objects and can be used as a robust feature descriptor for object reconstruction and classification. PG-Net adopts multi-task learning by estimating surface and volumetric representations. Surface and volumetric representations of objects have complementary properties of three-dimensional objects. Combining these modalities is more informative than using one modality alone. Objects are combinations of functional parts and part geometry is essential to synthesize each part of objects; therefore, PG-Net employs a part identifier to learn part geometry. Additionally, we augmented a dataset by interpolating individual parts, which helps learning part geometry and finding local/global minima of PG-Net. To demonstrate the capability of learnt object representations of PG-Net. We performed object reconstruction and classification tasks on two standard-large-scale datasets. PG-Net outperformed the other state-of-the-art methods in object synthesis, classification, and reconstruction.

© 2019 Elsevier B.V. All rights reserved.

### 1. Introduction

Technical advancement of 3D printing, virtual reality and augmented reality has brought significant interest of modeling and understanding three-dimensional object to the computer vision and graphics communities. This increased interest has led to progress of three-dimensional object synthesis [1, 2, 3], reconstruction [4, 5], and classification [6]. Moreover, emergence of neural networks and creation of large-scale three-dimensional object datasets [7, 8] inspired researchers to conduct three-dimensional object representation learning and synthesis by using view-based projections [9, 10], polygon meshes [11, 12, 13], point clouds [1, 14], and voxelized three-dimensional objects in voxel grids [15, 16]. In this work, we propose Part Geometry Network (PG-Net) which synthesizes realistic objects as a conditional generative model.

Many works [37, 38, 39] have been showed that jointly solving multiple tasks, named multi-task learning [19], helps im-

proving generalizability of estimation models. From this motivation, we adopt multi-task learning to optimize PG-Net.

Surface and volumetric representations of objects contain unique features and are complementary. Surface representation imposes boundary and connectivity information of each part of an object, and volumetric representation determines interior geometry which is used for heat flow calculation [17, 18]. However, using these two representations for multi-task learning has not been well-examined for learning compact object representations and synthesizing objects.

We used these modalities for multi-task learning to enhance generalizability of the model by designing the model to estimate surface and volumetric representations with the encoder-decoder structure. Knowing both surface and volumetric representations is crucial for learning a perceptual set of attributes, such as the connectivity of parts and the details of local interior regions, hence reducing defects in synthesized objects. In this way, the model learns surface properties to learn the interior volumetric representations of objects and vice-versa [19].

Part geometry of objects is critical to learn object distribution with parametric models [20, 21], since objects are combinations

\*Corresponding author: Tel.: +1-765-430-9721;  
e-mail: [kim2030@purdue.edu](mailto:kim2030@purdue.edu) (Sangpil Kim)

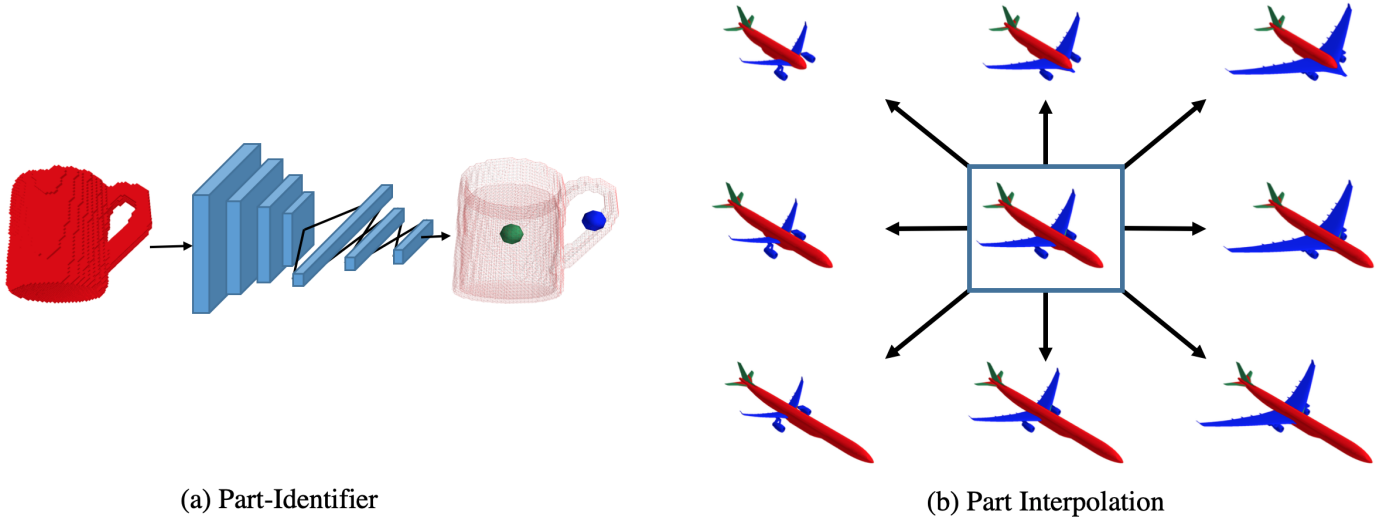


Fig. 1: Part-Identifier and part interpolation. (a) Part-Identifier is optimized to predict the location, volume, and surface area of parts. (b) Mesh models are expanded by relocating vertices of parts that expands searching space of optimization process.

of specific parts. However, learning part geometry is not trivial because defining the boundary of each part is a complex problem since parts can be connected in many different ways.

Each part has their own unique shape and location for their purpose, which creates unique part geometry. Therefore, we propose Part-Identifier which learns part geometry for predicting each part geometry (see Figure 1 a). Part-Identifier shares the part geometry information with other networks through back-propagating the meaningful gradients during training. With the understanding of part geometry, PG-Net can generate more realistic objects.

However, current datasets [8, 7] were assembled by collecting individual objects, and, therefore, each part of objects is unique. From this reason, the datasets are not sharing the information of parts among similar models, effectively. This makes it hard to learn relationship and connectivity between each part of objects. To alleviate the shortage of current datasets, we expanded the dataset [22] which has part labeled triangle mesh models by reshaping specific parts of objects as shown in Figure 1 (b).

Our main contributions are summarized as follows:

1. We improved generalizability of PG-Net for jointly estimating surface and volumetric representations as multi-task learning. In this way PG-Net learns complementary aspects of surface and volumetric properties of the object.
2. We propose learning method of part geometry, which improves the fidelity of each part of synthesized objects. Learning part geometry is done explicitly by optimizing the Part-Identifier which estimates each part geometry of objects.
3. We augmented a three-dimensional object dataset to expand the learning space of PG-Net by adjusting and scaling functional parts of objects. With expanded dataset, the optimizer can find better local/global minima of PG-Net.

Toward this end, we demonstrated ablation studies and comparison experiments with other methods. We performed an object synthesis with a one-hot encoded class vector and a vector from normal distribution. For evaluating shape representations obtained from PG-Net, we performed two applications: object reconstruction and classification. From our experiments, proposed method outperformed the other state-of-the-art methods in three-dimensional object synthesis, classification, and reconstruction.

## 2. Related Works

Generative models have been widely studied in the field of computer vision. There are two types of well-known generative models: Variational Auto-Encoder (VAE) [23] and Generative Adversarial Network (GAN) [24]. VAE consists of an encoder and decoder. The encoder encodes input information into a low dimensional vector which is perturbed by the Gaussian noise. In GAN, the model is optimized to mimic a data distribution by playing minimax game to find a Nash equilibrium [25] between the generator and discriminator; the generator is optimized to fool the discriminator by generating realistic data and the discriminator is optimized to identify faked data from the generator.

Generative models with object priors: generative models which synthesize 3D objects with object priors use an encoder-decoder structure for generating objects. Achlioptas *et al.* [1] proposed the encoder-decoder structure model, which learns point cloud representation by minimizing Earth Movers Distance and Chamfer Distance with the supervision of objects' class information. However, this method cannot synthesize 3D objects from a noise distribution, and must have reference models. Umetani *et al.* [26] synthesized deformed quad meshes from reference objects by exploring the manifold of the parameterized mesh surfaces with an encoder-decoder framework.

Liu *et al.* [27] proposed a method that reconstructs user-inputted 3D objects by mapping them into the hidden latent space and decoding it. Xie *et al.* [2] introduced an energy-based model which approximates the 3D shape probability distribution with Markov Chain Monte Carlo methods such as Langevin dynamics. Groueix *et al.* [3] generated the surface of 3D shapes with a generative model from point cloud objects or images. Kalogerakis *et al.* [28] and Carlson *et al.* [29] synthesized new 3D objects by reassembling the parts which are retrieved from an existing database with non-parametric approaches. The works from this section require reference objects to synthesize or deform objects. However, our proposed model synthesizes objects without observing objects or images as prior information.

Generative models without object priors: generative models without object priors use a decoder structure for generating synthetic objects. Wu *et al.* [30] proposed a generative adversarial loss with 3D volumetric convolution and synthesized novel 3D objects from the normal distribution. Smith *et al.* [31] improved 3D-GAN with Wasserstein GAN [32], which enhances the stability of the learning process. The conditional generative adversarial network (CGAN) [33] generates targeted images given a one-hot encoded class vector and noise vectors. Chen *et al.* [16] generated colored 3D objects in voxel grids given shape descriptions by jointly learning representations of the text description and 3D colored shapes with metric learning. Conditional Variational Auto-Encoder (CVAE) [34] has been used to learn specific patterns which are structured in the underlying data distribution. Bao *et al.* [35] combined the CVAE framework and adversarial loss, which performed fine-grained image generation. The conditional generative model for 3D objects has not been well-explored, and thus, existing methods in the 3D domain have no explicit controllability to sample a specific class of 3D objects with a single pipeline. In our method, a single pipeline is used to generate targeted objects given a targeted class one-hot vector.

Object reconstruction: various methods have been proposed for object reconstruction in three-dimensional domain. Dai *et al.* [4] reconstructed a corrupted distance field and state voxel grid with an encoder-decoder model and a shape database. The final output of this method was the reconstructed distance fields of 3D objects in the 3D space. Sung *et al.* [36] used the structure of 3D objects as prior knowledge for shape completion given noisy depth scans of objects. These works were specifically dedicated to solve shape reconstruction task.

### 3. Preliminary

Intersected Surface Area (ISA), which calculates intersected area within a cubic voxel for all voxels in defined voxel grids, was introduced by Yarotsky [12] and formulates as the following:

$$\int_{\Sigma \cap V_c} dS \quad (1)$$

, where  $V_c$  is a  $1 \times 1 \times 1$  cubic voxel, and  $\Sigma$  is the surface of an object. ISA contains the surface area value of each voxel grid

(see Figure 2 a). The Mean Curvature (MC), noted as  $H$ , is the mean value of maximum curvature  $\kappa_1$  and minimum curvature  $\kappa_2$ , which are extrinsic measures of curvature. We also tested 2-ring and 3-ring neighborhoods of vertices for mean curvature calculation in triangle mesh models, but the result was the best when using 1-ring neighborhood. Therefore, we used 1-ring neighborhood for mean curvature modality (see Figure 2 b).

Interior Volume (IV) is a collection of Boolean values in voxel grids (see Figure 2 c). In IV, if the cubic voxel is enclosed by the surface of an object; then, the voxel value is assigned as true. These three modalities have complimentary properties, and it is listed in Table 1.

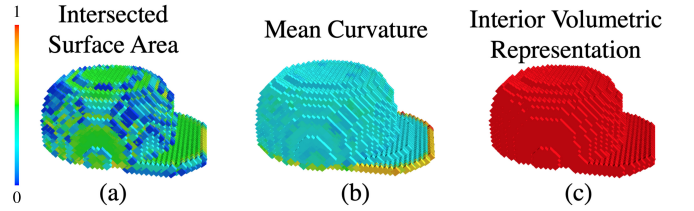


Fig. 2: Object representation modalities. (a) is the Intersected Surface Area (ISA) in a cubic voxel, (b) is the Mean Curvature (MC), and (c) is the Interior Volumetric (IV) representations of an object.

Table 1: Complementary properties of three different modalities.

Properties	IV	MC	ISA
Surface area information			✓
Association of vertices		✓	
Interior information	✓		

For notations in this paper,  $\mathcal{D}$ ,  $D$ , and  $R$ , are the discriminator, decoder, and refiner respectively. Generator is the combination of the decoder and refiner.  $H$  and  $isa$  are mean curvature and intersected surface area in a voxel, respectively.  $pl_i$ ,  $pv_i$  and  $ps_i$  are a  $\{x, y, z\}$  Cartesian coordinate of a central location, volume and surface area, respectively, where the lower index  $i \in \{\text{Body, Wheel, ... Legs}\}$  is the index of each part of objects. The  $\mathbf{z} \in \mathbb{R}^{200}$  is a vector of  $\mu + r \odot \exp(\epsilon)$ , where  $r \sim N(0, 1)$ ,  $N$  is normal distribution, and  $\odot$  is an element-wise multiplication.  $\mu$  and  $\epsilon$  are the mean and covariance of the latent vector from the encoder.  $p_d$  is data distribution and  $p_z$  is a prior distribution of the decoder.  $\mathbf{v}$  is the Boolean voxels from the true data distribution.

### 4. Multi-task Learning

In PG-Net, the decoder estimates three modalities, IV, MC, and ISA as multiple tasks. PG-Net estimates MC and ISA which are surface modalities because recent works [40, 41] showed that surface properties are informative for data-driven representation learning. From our experiment, unlike a volumetric-representation-alone framework, consolidating surface knowledge penalizes inaccurate surface estimates. Therefore, we designed the decoder to estimate MC and ISA of ob-

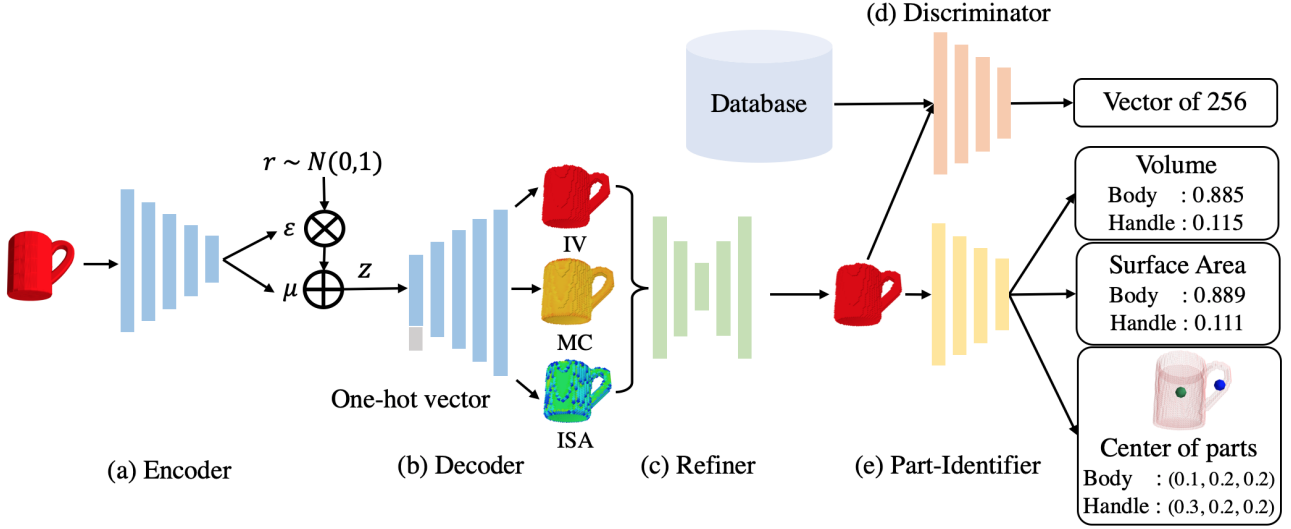


Fig. 3: Overview of PG-Net. The encoder encodes IV of the objects, and the decoder synthesizes IV, MC, and ISA. A diverse 3D object from a specified category is synthesized given a one-hot encoded class vector and a noise vector from the normal distribution. During the test stage, the encoder is removed from the pipeline. Noise vector  $\mathbf{z}$  is concatenated with a one-hot vector.

jects. We also experimented with Gaussian curvature and surface normal vectors, but the result was not better than using MC and ISA. As one of the tasks, the decoder was optimized to learn the surface representation by minimizing cost function  $\mathcal{L}_{surf}$ :

$$\sum \|\widehat{isa} - isa\|_1 + \lambda \cdot \sum \|\widehat{H} - H\|_1 \quad (2)$$

, where  $\lambda$  is a hyperparameter.

Another task was estimating volumetric representation of objects. PG-Net learns interior volume information of objects by minimizing sum of Sigmoid Cross-Entropy loss  $\mathcal{L}_{vol}$ :

$$-\mathbf{v} \log(\mathcal{D}(\mathbf{z})) - (1 - \mathbf{v}) \log(1 - \mathcal{D}(\mathbf{z})) \quad (3)$$

In terms of learning strategy, synthesizing IV, MC, and ISA with a single decoder is solving multiple tasks. In mathematical expression, PG-Net is optimized to jointly minimize the two cost functions:

$$\mathcal{L}_{vol} + \kappa \cdot \mathcal{L}_{surf} \quad (4)$$

, where  $\kappa$  is a hyperparameter. Jointly, learning both surface and volumetric representations is an informative way to discriminate objects because they provide complementary information of objects. Therefore, PG-Net uses the knowledge of surface and volumetric representations for learning object distribution in three-dimensional space.

## 5. Refiner and Discriminator

Adversarial training has been remarkably successful for synthesizing images and objects [42, 43, 44]. The goal is finding equilibrium between a generator and discriminator by playing minimax game between the generator and discriminator.

From this motivation, we experimented with an adversarial loss term [24] and found that high fidelity objects are generated. To enhance stability of adversarial training, we used least square adversarial loss [45]  $\mathcal{L}_{gan}$ :

$$\mathbb{E}_{\mathbf{z} \sim p_z}[(\mathcal{D}(R(D(\mathbf{z}))) - 1)^2] + \mathbb{E}_{\mathbf{v} \sim p_d}[(D(\mathbf{v}) - 1)^2] + \mathbb{E}_{\mathbf{z} \sim p_z}[(\mathcal{D}(R(D(\mathbf{z}))))^2] \quad (5)$$

Discriminator discriminates fake objects from model distribution and true objects from database. To train the discriminator, we minimized the cost function:

$$\mathbb{E}_{\mathbf{v} \sim p_d}[(D(\mathbf{v}) - 1)^2] + \mathbb{E}_{\mathbf{z} \sim p_z}[(\mathcal{D}(R(D(\mathbf{z}))))^2] \quad (6)$$

To synthesize objects given  $\mathbf{z}$ , we minimize the gap between  $p(\mathbf{z})$  and three-dimensional object distribution in the latent space [35] by minimizing  $\mathcal{L}_{KL}$ :

$$\mu^T \mu + \text{sum}(\exp(\epsilon) - \epsilon - 1) \quad (7)$$

Refiner is composed of an hourglass architecture [46] and refines coarse objects from the decoder by consolidating estimated three modalities. We found that sometimes estimation of three modalities are not accurate nor consistent. These problems were solved in a more intelligent way. We designed the refiner to penalize poorly estimated modalities from the decoder with the loss term  $\mathcal{L}_{ref}$ :

$$-\mathbf{v} \log(R(D(\mathbf{z}))) - (1 - \mathbf{v}) \log(1 - R(D(\mathbf{z}))) \quad (8)$$

For training the refiner, we jointly trained it with Part-Identifier by optimizing the cost function:

$$\alpha \cdot \mathcal{L}_{gan} + \beta \cdot \mathcal{L}_p + \gamma \cdot \mathcal{L}_{ref} + \mathcal{L}_{KL} \quad (9)$$



## 6. Learning Part Geometry

Part geometry is an important factor to understand three-dimensional object space [36]. From this motivation, we propose Part-Identifier for learning part geometry. To expand the learning space of each part of object, we further augmented the objects by interpolating the core parts of objects.

### 6.1. Part-Identifier

Part-Identifier in PG-Net estimates part geometry information (see Figure 3 e), and the information is back-propagated through the pipeline of PG-Net. For learning the part geometry, we considered using point-wise segmentation of parts. However, it was computationally expensive, and therefore, we regressed the center location, surface area, and the volume of each part. The position of the central coordinate, surface area, and the volume of each part are normalized by dividing the resolution of voxel grids, the surface area, and volume of the object, respectively. Volume and surface area of parts impose meaningful information of objects such as relations of each part. Part-Identifier implicitly learns the relationship between different parts and part geometry by estimating the locations, volumes and surface areas of parts. This module guides the generator towards lower local/global minima by minimizing the  $\mathcal{L}_p$ :

$$\sum_i \|\widehat{pl}_i - pl_i\|_1 + \|\widehat{ps}_i - ps_i\|_1 + \|\widehat{pv}_i - pv_i\|_1 \quad (10)$$

,where  $i \in \{\text{Body, Wheel, ... Legs}\}$ , the index of each part.

Table 2: Parts for *Expanded* dataset creation.

Classes	Interpolated Parts		
Airplane	Body	Wings	
Bag	Handle	Case	
Cap	Crown	Brim	
Car	Roof	Wheel	Hood
Chair	Back	Seat	Leg
Lamp	Legs	Base	Shade
Mug	Cup	Handle	
Table	Top	Leg	

### 6.2. Part expansion

The goal of part expansion as shown in Figure 4 is to capture part geometry effectively and expand search space for optimizing PG-Net. In order to create our dataset, named *Expanded*, we expanded Projective [22] which is a subset of the ShapeNet [7]. Projective consists of mesh models, and each part of the object is vertice-wise labeled. We first chose parts (see Table 2) from each class as listed in the Table 3. We removed outliers of the labeled vertices from Projective with Density-Based Spatial Clustering of Applications with Noise (DBSCAN). Additionally, we manually filtered out outliers in the dataset. Then, we adjusted the triangle mesh vertices to augment the dataset. After adjusting the triangle mesh vertices, we manually filtered out odd-looking objects. For example, we removed an object

where the table top surface area is smaller than the surface area of the leg. To complete our dataset expansion, we added the models from the Scalable dataset [22], which is a subset of the ShapeNet into *Expanded* to increase diversity of objects. The statistics of our dataset are detailed in Table 3.

Table 3: Statistics of *Expanded*, Scalable [47], and Projective [22].

Classes	Expanded (Ours)	Scalable [47]	Projective [22]
Airplane	36,018	2,690	500
Bag	666	76	76
Cap	477	55	55
Car	12,739	878	500
Chair	13,041	3,746	500
Lamp	13,014	1,546	500
Mug	1,629	184	184
Table	30,015	5,263	500
Total	107,599	14,438	2,815

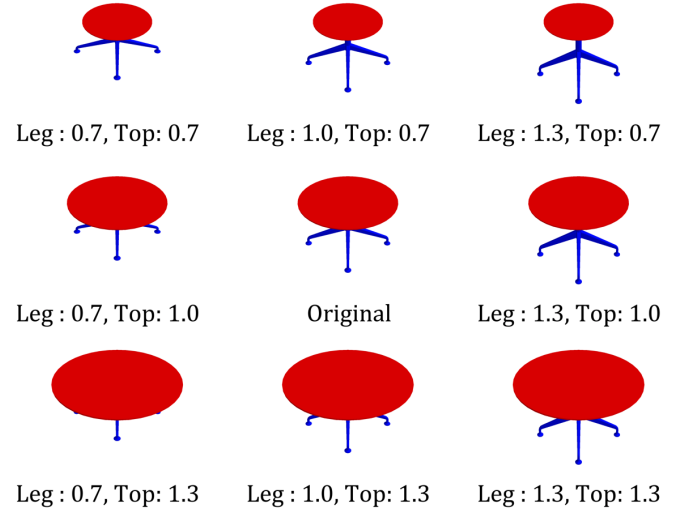


Fig. 4: Example of part expansion. Each part of tables is expanded by relocating vertices of parts to share each part information within the object. Parts are numbered based on their scaling proportion.

## 7. Experiments

In this section, we present experimental validations and analysis of three-dimensional object synthesis, classification, and reconstruction to validate the efficacy of PG-Net. We used two standard large-scale three-dimensional object dataset: ModelNet [8] and *Expanded* (see Table 3), which is an expanded version of Projective [22]. For dataset splitting, we divided our dataset into three sets: 70% as a training set, 10% as a validation set, and 20% as a test set. We performed isosurface and Laplacian smoothing for visualizing the objects and used  $L_1$  metric for quantitative evaluations.

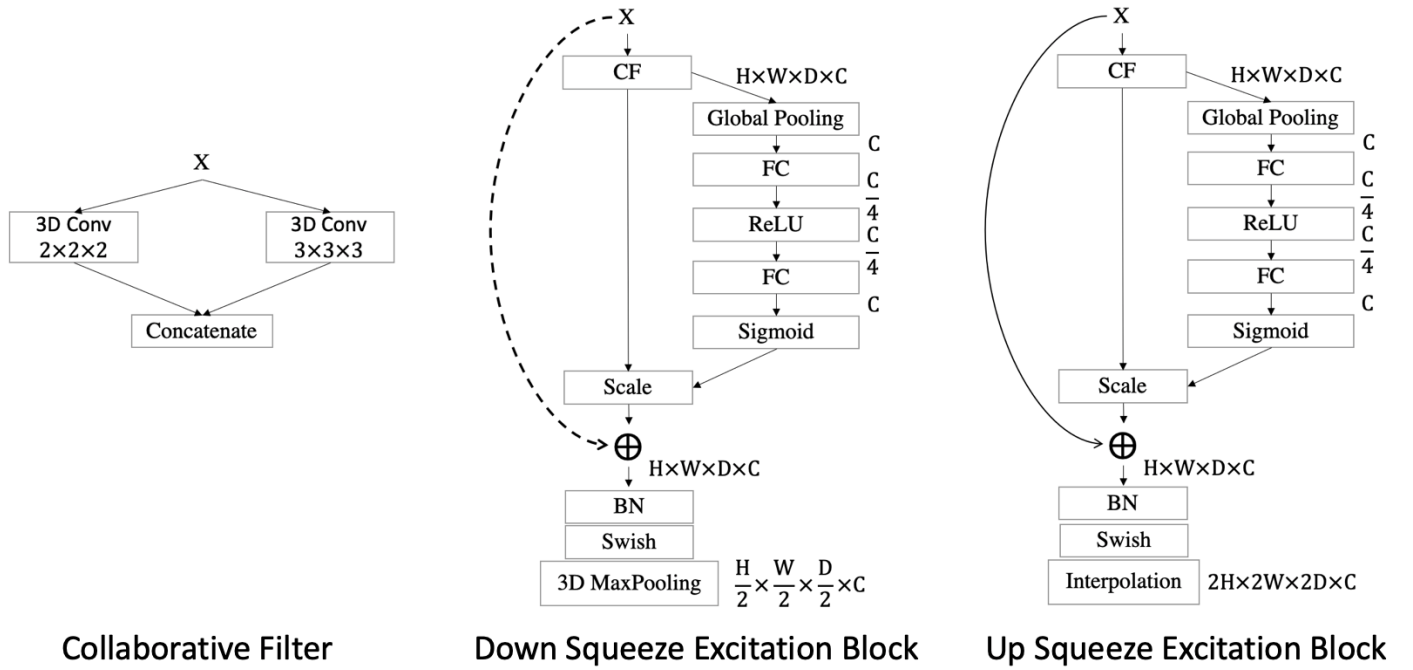


Fig. 5: Block diagrams of Collaborative Filters (CF) and Squeeze-and-Excitation block [48]. The dotted shortcuts increase dimensions if channels of  $X$  is smaller than  $C$ . BN is Batch Normalization, and  $C$  is the number of channels. Kernel size of  $1 \times 1 \times 1$  3D convolution is applied to the solid shortcut when there are mismatches in the number of channels between  $X$  and  $C$ .

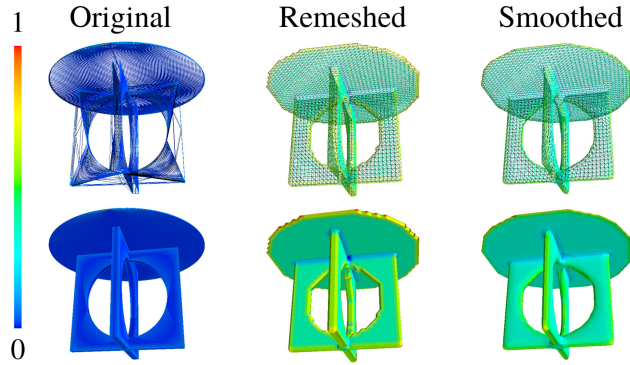


Fig. 6: Comparing an original mesh with an isosurface from a voxelized object. The first column shows mean curvatures from an original mesh. The second column shows mean curvatures from the mesh generated by marching-cubes algorithm given a voxelized object in  $40^3$  grids. The third column shows mean curvatures from the smoothed mesh by applying Laplacian smoothing with the second column mesh.

### 7.1. Dataset preprocessing

We preprocessed mesh models to ensure the existence of mean curvature per each voxel grid. First, we voxelized mesh models with the resolution of  $40^3$  voxel grids, and then we used the marching-cubes algorithm [49] to regenerate triangle meshes since voxelized objects are deformed from the original shape (see Figure 6 Remeshed). We further smoothed the meshes with Laplacian smoothing [50] to reduce the variation of mean curvature (see Figure 6 Smoothed). From the preprocessed meshes, we calculated MC and ISA and assigned in  $40^3$

voxel grids. For the ISA extraction, we used SurfaceArea operation in a library [12] that computes face areas within each cubic voxel. We obtained MC by calculating the average of the minimum and maximum principal curvatures with 1-ring neighborhood distance. After obtaining all data, we normalized and scaled them to be within the interval  $[-1, 1]$ .

### 7.2. Architecture details

In neural network, salient parts are critical aspects that differentiate objects from each other. However, salient parts can have large variation spatially in three-dimensional space and are dependent of the kernel sizes of convolutional layers [51]. Because of this variation, choosing a right kernel size is an important factor to extract robust features from the input data. Therefore, state-of-the-art image classifiers [52, 48] use filters with multiple sizes operating on the same level. From this motivation, we used Collaborate Filter (CF) as shown in Figure 5 to capture salient parts effectively from diverse regions of the input voxels.

Convolution neural network fuses channel-wise information within local receptive fields at each layer to construct robust features, and squeeze and excitation block [48] adaptively adjusts channel-wise feature responses. Therefore, we used Up Squeeze Excitation Block (USEB) and Down Squeeze Excitation Block (DSEB) (see Figure 5). Definitions of network architecture are defined in Table 4. Throughout the network, zero padding was applied when the size of the output from a previous layer is not divisible by two. We used batch normalization

layer and Swish [53] activation function in between the transition layers. In discriminator, Swish activation function was replaced with the LeakyReLU activation function on each DSEB layer.

Table 4: Definitions of network architecture. A, B, C, D, and E are the encoder, decoder, refiner, discriminator, and Part-Identifier, respectively. In the type column, BI, FC, and COV, are abbreviations of bilinear interpolation, fully-connected layer, and convolutional layer, respectively.

	Type	Filter/Stride	Output Size	#Channel
A	Input		40x40x40	1
	DSEB	CF/1	20x20x20	16
	DSEB	CF/1	10x10x10	32
	DSEB	CF/1	5x5x5	64
	DSEB	CF/1	3x3x3	128
	3D-COV	1x1x1/1	2x2x2	256
	3D-COV	1x1x1/1	1x1x1	512
B	2D-COV	1x1/2	1x1	300
	BI		2x2	300
	2D-COV	3x3/2	2x2	150
	BI		5x5	150
	2D-COV	3x3/2	2x2	40
	BI		10x10	40
	Reshape		5x5x5	32
	USEB	CF/1	10x10x10	16
	USEB	CF/1	20x20x20	8
	USEB	CF/1	40x40x40	3
C	3D-COV	3x3x3/2	20x20x20	8
	3D-COV	3x3x3/2	10x10x10	16
	3D-COV	3x3x3/2	5x5x5	32
	3D-COV	3x3x3/1	5x5x5	16
	BI		10x10x10	16
	3D-COV	3x3x3/1	10x10x10	8
	BI		20x20x20	8
	3D-COV	3x3x3/1	20x20x20	1
	BI		40x40x40	1
D/E	DSEB	CF/1	20x20x20	64
	DSEB	CF/1	10x10x10	128
	DSEB	CF/1	5x5x5	256
	DSEB	CF/1	3x3x3	512
	FC		512	
	FC		256(D)/95(E)	

### 7.3. Implementation

We used two NVIDIA TITAN Xp GPUs and an Intel i7-6850K CPU with 64GB of RAM for all of our experiments. The artificial neural network was developed with TensorFlow deep learning framework [54] which was accelerated by the CUDA instruction for the GPU computation. The networks were optimized by using the ADAM optimizer [55] with the initial parameters: learning rate = 0.0025,  $\beta_1 = 0.9$ , and  $\beta_2 = 0.999$ . For the hyperparameters, we used  $\alpha = 0.5$ ,  $\beta = 0.1$ ,  $\lambda = 10^{-4}$ ,  $\kappa = 1$ , and  $\gamma = 0.1$ .

We trained PG-Net in two stages with a batch size of 16. In the first stage, we trained the encoder and decoder separately

from other networks, which were converged approximately after 250 epochs. Then we stacked the refiner, discriminator, and Part-Identifier for the second-stage training, which requires approximately 200 epochs to converge. During the second-stage training, we did not update the encoder and decoder. In order to improve the stability of training the discriminator, we updated the refiner and Part-Identifier twice per each discriminator update to enhance the stability of the learning process [25]. After training the second-stage, we jointly updated all networks with the learning rate of  $10^{-7}$ , which required approximately 100 epochs to converge. We used the initial learning rates of 0.0005 and 0.0025 for the first and second-stage training, respectively. We dropped the running rate by half in every 25 epochs. For synthesizing objects, we synthesized objects from  $\mathbf{z} \in \mathbb{R}^{200}$  which was sampled from a normal distribution  $N(0, 1)$  and a one-hot encoded class vector.

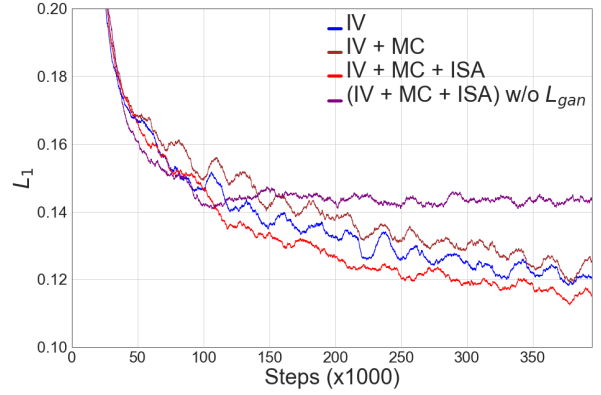


Fig. 7:  $L_1$  loss plots of four experimental models on the test set.

Table 5: Quantitative results as  $L_1$  metric of baselines and other methods for 3D object generation without 3D object references. Lower value is better.

	Seen	Unseen
3D-CIWGAN [31]	0.225	0.232
3D-CVAE [56]	0.173	0.199
PG-Net w/o Part-Identifier	0.164	0.194
PG-Net w/o Part Interpolation	0.183	0.204
PG-Net w/o Refiner	0.112	0.135
PG-Net (Ours)	<b>0.083</b>	<b>0.117</b>

### 7.4. Ablation study

**Does multi-task and adversarial learning lead to the synthesis of better quality objects?** For the ablation study, we ran four experiments: (i) IV, which only estimates volumetric representations of objects for training; (ii) IV+MC, which estimates IV and MC; (iii) IV+MC+ISA, which estimates IV, MC and ISA; (iv) (IV+MC+ISA) w/o  $L_{gan}$ , which is the same as (iii) but without  $L_{gan}$ . Figure 7 shows the performances of all experiments. The results indicate that the network learns a structural correlation across the local surface and volume descriptors to

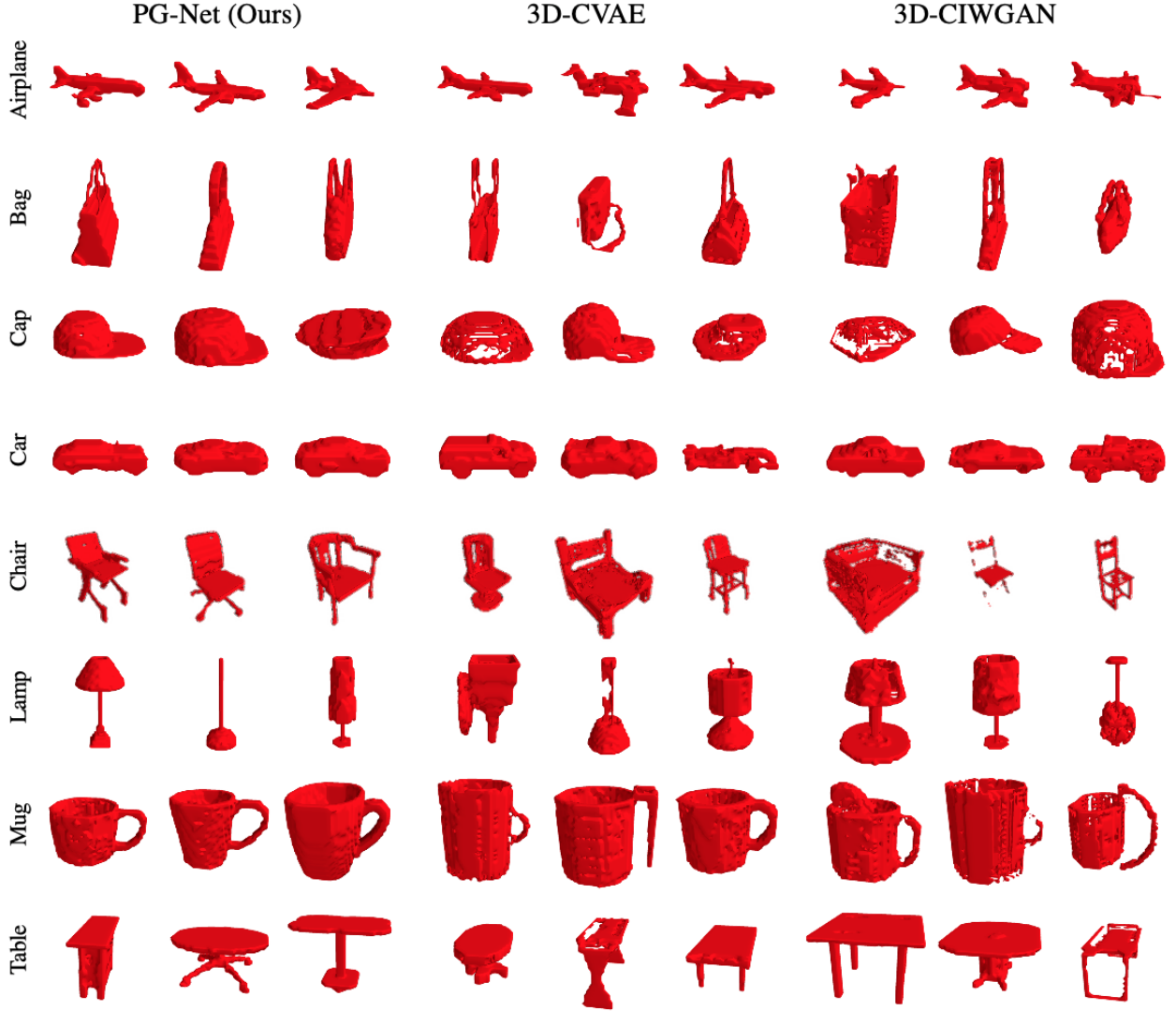


Fig. 8: Objects generated from the model given a one-hot encoded class vector and a vector from the normal distribution without a reference object. The resolution of the voxel grids was  $40^3$ , and the objects were binarized by the threshold of 0.5.



Fig. 9: Sampled objects from the models w/ Part-Identifier and w/o Part-Identifier given normal distribution and a one-hot encoded class vector.

improve the fidelity of final outputs. This ablation study validates the rationale of multi-task and adversarial learning with IV, MC, and ISA modalities.

**Why learn the geometry of parts?** Furthermore, we explored the efficacy of learning part geometry (see Figure 9). Table 5 shows the quantitative evaluation of the proposed method for the following baselines: (i) *w/o Part-Identifier*, which does not have the part identifier; (ii) *w/o Part interpolation*, which is the same pipeline as PG-Net but uses dataset without data

augmentation with part interpolating; (iii) *w/o Refiner*, which is trained without the refiner to evaluate the efficacy of the refiner. From the result shown in Table 5, we verified that *part identifier* improves the quality of synthesized objects since  $PG-Net$  of  $L_1$  are lower than *w/o Part-Identifier*. Also, PG-Net gave lower metric scores than *w/o Part interpolation* and *w/o Refiner*, which directly shows that the refiner and part interpolation results in learning robust shape representations. As a conclusion, PG-Net performs better than the other baselines, and each of our proposed methods reduces artifacts and holes of synthesized objects.

### 7.5. Synthesizing 3D objects

We conditionally generate 3D objects by sampling a latent vector  $\mathbf{z}$ , which was mapped to the object space, and a one-hot encoded class vector. We compared our method with 3D-CVAE [56] and 3D-CIWGAN [31]. For 3D-CVAE, we followed CVAE base model in [56] and used the same



encoder/decoder definitions as PG-Net. We combined 3D-IWGAN [31] and CGAN [33] for 3D-CIWGAN. All methods used *Expanded* with an one-hot encoded class vector as a condition. Figure 8 shows the synthesized objects from our method and the others. Since we were sampling objects from noise distribution without ground truth, displaying the exact same objects for each method was impossible. Each voxel value of the synthesized objects was binarized with a threshold of 0.5. We visualized volumetric data after applying the marching-cubes algorithm and Laplacian smoothing. We observed fewer artifacts and holes in the objects from PG-Net than the objects from the other methods. 3D-CIWGAN is not able to synthesize the nose of the airplane, and 3D-CVAE was effective on airplanes but had many holes in the table class. Unlike 3D-CVAE and 3D-CIWGAN, PG-Net preserved part geometry and well-defined surfaces.

### 7.6. 3D Object classification

We evaluated object representations from PG-Net by performing 3D object classification on the ModelNet [8] which offers two kinds of datasets: ModelNet10 and ModelNet40. ModelNet10 and ModelNet40 comprises 4,899 objects and 10 classes and 12,331 objects and 40 classes, respectively. For unsupervised training, we used the same method and a dataset for fine-tuning PG-Net from Achlioptas *et al.* [1]. The dataset consists of 57,000 objects from 55 categories in ShapeNet [7]. We fine-tuned the optimized PG-Net without Part-Identifier since the parts' labels are not available in the dataset. After fine-tuning PG-Net, we extracted features from the last two layers of encoder and concatenated them to create high-dimensional representations. Then we classified the representations with a linear SVM trained on the 3D classification benchmark [8]. Our method outperforms other existing works as shown in Table 6. From the results, PG-Net extracts robust features which can effectively distinguish objects, and, therefore, PG-Net can be used as a feature descriptor for the object analysis and other applications.

Table 6: The comparison on classification accuracy between PG-Net and the other unsupervised methods. Higher number is better.

	ModelNet10	ModelNet40
SPH [57]	79.8%	68.2
VConv-DAE [58]	80.5%	75.5%
3D-GAN [30]	91.0%	83.3%
GMPC [1]	95.4%	84.5%
ECC [59]	90.0%	83.2%
FoldingNet [60]	94.4%	88.4%
PG-Net (Ours)	<b>95.6%</b>	<b>89.1%</b>

### 7.7. Object reconstruction

We performed three-dimensional object reconstruction to evaluate the efficacy of PG-Net and compared the results with 3D-EPN [4]. For our experiment, we used a corrupted dataset from 3D-EPN. We interpolated the input objects from  $32^3$  voxel

grids into  $40^3$  voxel grids to match the input resolution of PG-Net. For the quantitative evaluation, we converted the reconstructed objects into Boolean voxels with the threshold of 0.5. Then we counted wrongly estimated voxels and divided with the total number of occupied voxels in ground truth. We used the pre-trained weights of EPN-unet w/ class version from the 3D-EPN project page as a comparison. The quantitative results of PG-Net and 3D-EPN are compared in Table 7. The error values of PG-Net are lower than those of 3D-EPN. As shown in Figure 10 PG-Net also shows better qualitative results as compared to 3D-EPN. From quantitative and qualitative experiment results, PG-Net outperformed 3D-EPN in large margin along all classes. Therefore, PG-Net, which was trained with multi-task and part geometry learning method, effectively learns object distribution and nicely reconstructs three-dimensional corrupted object.

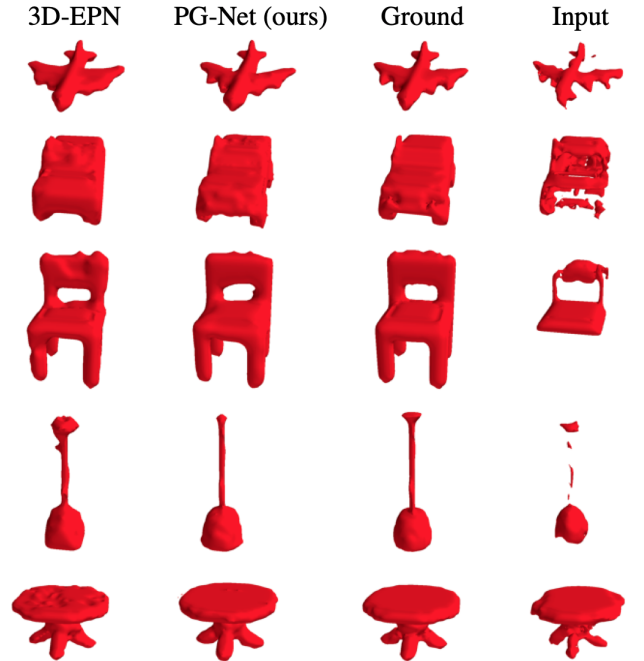


Fig. 10: Qualitative results of three-dimensional object reconstruction experiment between PG-Net and 3D-EPN.

Table 7: Quantitative results of reconstruction task and comparison of 3D-EPN [4] and PG-Net. Lower value is better.

Class (# of train / # of test )	3D-EPN [4]	PG-Net (Ours)
Air. (3.3K / 0.8K)	0.226	<b>0.202</b>
Car (5K / 1K)	0.197	<b>0.191</b>
Chair (5K / 1K)	0.309	<b>0.273</b>
Lamp (1.8K / 0.5K)	0.407	<b>0.392</b>
Table (5K / 1K)	0.338	<b>0.251</b>
Total (20.1K / 4.3K)	0.286	<b>0.249</b>

## 8. Conclusions and Discussion

In this paper we propose PG-Net which was optimized with multi-task learning and part geometry learning for object synthesis. Results from our study suggest that multi-task learning increases the fidelity of generated objects and learning part geometry enhances the realism of each part of the synthesized objects. PG-Net exceeded the other state-of-the-art methods in object synthesis, classification, and reconstruction. As limitations, 3D convolution layers with voxels are computationally expensive and, thus, require ample memory. However, these limitations can be solved by using a recurrent neural network with polygon mesh vertices or point clouds representations. For a future work, fusing other surface modalities as an input with point clouds representation could be further explored with a recurrent neural network.

## References

- [1] Achlioptas, P, Diamanti, O, Mitliagkas, I, Guibas, L. Learning representations and generative models for 3d point clouds. In: International Conference on Machine Learning. 2018, p. 40–49.
- [2] Xie, J, Zheng, Z, Gao, R, Wang, W, Zhu, SC, Wu, YN. Learning descriptor networks for 3d shape synthesis and analysis. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2018, p. 8629–8638.
- [3] Groueix, T, Fisher, M, Kim, VG, Russell, BC, Aubry, M. Atlasnet: A papier-m<sup>ch</sup> approach to learning 3d surface generation. arXiv preprint arXiv:180205384 2018;.
- [4] Dai, A, Ruizhongtai Qi, C, Nießner, M. Shape completion using 3d-encoder-predictor cnns and shape synthesis. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2017, p. 5868–5877.
- [5] Choy, CB, Xu, D, Gwak, J, Chen, K, Savarese, S. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In: European conference on computer vision. Springer; 2016, p. 628–644.
- [6] Chang, AX, Funkhouser, T, Guibas, L, Hanrahan, P, Huang, Q, Li, Z, et al. Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:151203012 2015;.
- [7] Chang, AX, Funkhouser, T, Guibas, L, Hanrahan, P, Huang, Q, Li, Z, et al. ShapeNet: An Information-Rich 3D Model Repository. Tech. Rep. arXiv:1512.03012 [cs.GR]; Stanford University — Princeton University — Toyota Technological Institute at Chicago; 2015.
- [8] Wu, Z, Song, S, Khosla, A, Yu, F, Zhang, L, Tang, X, et al. 3d shapenets: A deep representation for volumetric shapes. In: Proceedings of the IEEE conference on computer vision and pattern recognition. 2015, p. 1912–1920.
- [9] Kanezaki, A, Matsushita, Y, Nishida, Y. Rotationnet: Joint object categorization and pose estimation using multiviews from unsupervised viewpoints. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2018, p. 5010–5019.
- [10] Wu, Z, Zhang, Y, Zeng, M, Qin, F, Wang, Y. Joint analysis of shapes and images via deep domain adaptation. Computers & Graphics 2018;70:140–147.
- [11] Luciano, L, Hamza, AB. A global geometric framework for 3d shape retrieval using deep learning. Computers & Graphics 2019;79:14–23.
- [12] Yarotsky, D. Geometric features for voxel-based surface recognition. arXiv preprint arXiv:170104249 2017;.
- [13] Masci, J, Boscaini, D, Bronstein, M, Vandergheynst, P. Geodesic convolutional neural networks on riemannian manifolds. In: Proceedings of the IEEE international conference on computer vision workshops. 2015, p. 37–45.
- [14] Qi, CR, Yi, L, Su, H, Guibas, LJ. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in Neural Information Processing Systems. 2017, p. 5099–5108.
- [15] Girdhar, R, Fouhey, DF, Rodriguez, M, Gupta, A. Learning a predictable and generative vector representation for objects. In: European Conference on Computer Vision. Springer; 2016, p. 484–499.
- [16] Chen, K, Choy, CB, Savva, M, Chang, AX, Funkhouser, T, Savarese, S. Text2shape: Generating shapes from natural language by learning joint embeddings. arXiv preprint arXiv:180308495 2018;.
- [17] Kotte, A, Van Wieringen, N, Legendijk, J. Modelling tissue heating with ferromagnetic seeds. Physics in Medicine & Biology 1998;43(1):105.
- [18] Nelson, D, Charbonnel, S, Curran, A, Marttila, E, Fiala, D, Mason, P, et al. A high-resolution voxel model for predicting local tissue temperatures in humans subjected to warm and hot environments. Journal of Biomechanical Engineering 2009;131(4):041003.
- [19] Thrun, S, Pratt, L. Learning to learn. Springer Science & Business Media; 2012.
- [20] Bu, S, Han, P, Liu, Z, Han, J, Lin, H. Local deep feature learning framework for 3d shape. Computers & Graphics 2015;46:117–129.
- [21] Mo, K, Zhu, S, Chang, AX, Yi, L, Tripathi, S, Guibas, LJ, et al. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2019, p. 909–918.
- [22] Kalogerakis, E, Averkiou, M, Maji, S, Chaudhuri, S. 3d shape segmentation with projective convolutional networks. In: Proc. CVPR; vol. 1. 2017, p. 8.
- [23] Kingma, DP, Welling, M. Auto-encoding variational bayes. arXiv preprint arXiv:13126114 2013;.
- [24] Goodfellow, I, Pouget-Abadie, J, Mirza, M, Xu, B, Warde-Farley, D, Ozair, S, et al. Generative adversarial nets. In: Advances in neural information processing systems. 2014, p. 2672–2680.
- [25] Goodfellow, I. Nips 2016 tutorial: Generative adversarial networks. arXiv preprint arXiv:170100160 2016;.
- [26] Umetani, N. Exploring generative 3d shapes using autoencoder networks. In: SIGGRAPH Asia 2017 Technical Briefs. ACM; 2017, p. 24.
- [27] Liu, J, Yu, F, Funkhouser, T. Interactive 3d modeling with a generative adversarial network. arXiv preprint arXiv:170605170 2017;.
- [28] Kalogerakis, E, Chaudhuri, S, Koller, D, Koltun, V. A probabilistic model for component-based shape synthesis. ACM Transactions on Graphics (TOG) 2012;31(4):55.
- [29] Carlson, WE. An algorithm and data structure for 3d object synthesis using surface patch intersections. ACM SIGGRAPH Computer Graphics 1982;16(3):255–263.
- [30] Wu, J, Zhang, C, Xue, T, Freeman, B, Tenenbaum, J. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In: Advances in Neural Information Processing Systems. 2016, p. 82–90.
- [31] Smith, E, Meger, D. Improved adversarial systems for 3d object generation and reconstruction. arXiv preprint arXiv:170709557 2017;.
- [32] Arjovsky, M, Chintala, S, Bottou, L. Wasserstein gan. arXiv preprint arXiv:170107875 2017;.
- [33] Mirza, M, Osindero, S. Conditional generative adversarial nets. arXiv preprint arXiv:14111784 2014;.
- [34] Kingma, DP, Mohamed, S, Rezende, DJ, Welling, M. Semi-supervised learning with deep generative models. In: Advances in Neural Information Processing Systems. 2014, p. 3581–3589.
- [35] Bao, J, Chen, D, Wen, F, Li, H, Hua, G. Cvae-gan: fine-grained image generation through asymmetric training. CoRR, abs/170310155 2017;5.
- [36] Sung, M, Kim, VG, Angst, R, Guibas, L. Data-driven structural priors for shape completion. ACM Transactions on Graphics (TOG) 2015;34(6):175.
- [37] Li, S, Chan, AB. 3d human pose estimation from monocular images with deep convolutional neural network. In: Asian Conference on Computer Vision. Springer; 2014, p. 332–347.
- [38] Wu, Z, Valentini-Botinhao, C, Watts, O, King, S. Deep neural networks employing multi-task learning and stacked bottleneck features for speech synthesis. In: 2015 IEEE international conference on acoustics, speech and signal processing (ICASSP). IEEE; 2015, p. 4460–4464.
- [39] Zhang, C, Zhang, Z. Improving multiview face detection with multi-task deep convolutional neural networks. In: IEEE Winter Conference on Applications of Computer Vision. IEEE; 2014, p. 1036–1041.
- [40] Hanocka, R, Hertz, A, Fish, N, Giryas, R, Fleishman, S, Cohen-Or, D. Meshcnn: a network with an edge. ACM Transactions on Graphics (TOG) 2019;38(4):90.
- [41] Kostrikov, I, Jiang, Z, Panozzo, D, Zorin, D, Bruna, J. Surface networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2018, p. 2540–2548.
- [42] Cui, J, Li, S, Xia, Q, Hao, A, Qin, H. Learning multi-view manifold

- for single image based modeling. *Computers & Graphics* 2019;.
- [43] Park, E, Yang, J, Yumer, E, Ceylan, D, Berg, AC. Transformation-grounded image generation network for novel 3d view synthesis. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, p. 3500–3509.
  - [44] Liu, C, Yang, Z, Xu, F, Yong, JH. Image generation from bounding box-represented semantic labels. *Computers & Graphics* 2019;81:32–40.
  - [45] Mao, X, Li, Q, Xie, H, Lau, RY, Wang, Z, Smolley, SP. On the effectiveness of least squares generative adversarial networks. *arXiv preprint arXiv:171206391* 2017;.
  - [46] Newell, A, Yang, K, Deng, J. Stacked hourglass networks for human pose estimation. In: *European conference on computer vision*. Springer; 2016, p. 483–499.
  - [47] Yi, L, Kim, VG, Ceylan, D, Shen, I, Yan, M, Su, H, et al. A scalable active framework for region annotation in 3d shape collections. *ACM Transactions on Graphics (TOG)* 2016;35(6):210.
  - [48] Hu, J, Shen, L, Sun, G. Squeeze-and-excitation networks. *arXiv preprint arXiv:170901507* 2017;.
  - [49] Lewiner, T, Lopes, H, Vieira, AW, Tavares, G. Efficient implementation of marching cubes' cases with topological guarantees. *Journal of graphics tools* 2003;8(2):1–15.
  - [50] Vollmer, J, Mencl, R, Mueller, H. Improved laplacian smoothing of noisy surface meshes. In: *Computer graphics forum*; vol. 18. Wiley Online Library; 1999, p. 131–138.
  - [51] Paszke, A, Chaurasia, A, Kim, S, Culurciello, E. Enet: A deep neural network architecture for real-time semantic segmentation. *arXiv preprint arXiv:160602147* 2016;.
  - [52] Szegedy, C, Vanhoucke, V, Ioffe, S, Shlens, J, Wojna, Z. Rethinking the inception architecture for computer vision. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016, p. 2818–2826.
  - [53] Ramachandran, P, Zoph, B, Le, QV. Swish: a self-gated activation function. *arXiv preprint arXiv:171005941* 2017;.
  - [54] Abadi, M, Barham, P, Chen, J, Chen, Z, Davis, A, Dean, J, et al. Tensorflow: A system for large-scale machine learning. In: *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*. 2016, p. 265–283.
  - [55] Kingma, D, Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* 2014;.
  - [56] Yan, X, Yang, J, Sohn, K, Lee, H. Attribute2image: Conditional image generation from visual attributes. In: *European Conference on Computer Vision*. Springer; 2016, p. 776–791.
  - [57] Kazhdan, M, Funkhouser, T, Rusinkiewicz, S. Rotation invariant spherical harmonic representation of 3 d shape descriptors. In: *Symposium on geometry processing*; vol. 6. 2003, p. 156–164.
  - [58] Sharma, A, Grau, O, Fritz, M. Vconv-dae: Deep volumetric shape learning without object labels. In: *European Conference on Computer Vision*. Springer; 2016, p. 236–250.
  - [59] Simonovsky, M, Komodakis, N. Dynamic edgeconditioned filters in convolutional neural networks on graphs. In: *Proc. CVPR*. 2017;.
  - [60] Yang, Y, Feng, C, Shen, Y, Tian, D. Foldingnet: Point cloud auto-encoder via deep grid deformation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018, p. 206–215.