

BSP-Net: Generating Compact Meshes via Binary Space Partitioning

Zhiqin Chen
Simon Fraser University

Andrea Tagliasacchi
Google Research

Hao Zhang
Simon Fraser University

Abstract

Polygonal meshes are ubiquitous in the digital 3D domain, yet they have only played a minor role in the deep learning revolution. Leading methods for learning generative models of shapes rely on implicit functions, and generate meshes only after expensive iso-surfacing routines. To overcome these challenges, we are inspired by a classical spatial data structure from computer graphics, Binary Space Partitioning (BSP), to facilitate 3D learning. The core ingredient of BSP is an operation for recursive subdivision of space to obtain convex sets. By exploiting this property, we devise *BSP-Net*, a network that learns to represent a 3D shape via convex decomposition. Importantly, *BSP-Net* is unsupervised since no convex shape decompositions are needed for training. The network is trained to reconstruct a shape using a set of convexes obtained from a BSP-tree built on a set of planes. The convexes inferred by *BSP-Net* can be easily extracted to form a polygon mesh, without any need for iso-surfacing. The generated meshes are compact (i.e., low-poly) and well suited to represent sharp geometry; they are guaranteed to be watertight and can be easily parameterized. We also show that the reconstruction quality by *BSP-Net* is competitive with state-of-the-art methods while using much fewer primitives. Code is available at <https://github.com/czq142857/BSP-NET-original>.

1. Introduction

Recently, there has been an increasing interest in representation learning and generative modeling for 3D shapes. Up to now, deep neural networks for shape analysis and synthesis have been developed mainly for voxel grids [14, 18, 48, 50], point clouds [1, 33, 34, 55, 56], and implicit functions [5, 13, 22, 28, 52]. As the dominant 3D shape representation for modeling, display, and animation, polygonal meshes have not figured prominently amid these developments. One of the main reasons is that the non-uniformity and irregularity of triangle tessellations do not naturally support conventional convolution and pooling operations [19]. However, compared to voxels and point clouds, meshes can provide a more seamless and coherent

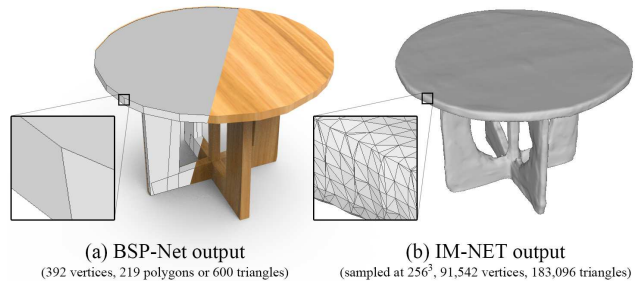


Figure 1: (a) 3D shape auto-encoding by *BSP-Net* quickly reconstructs a *compact*, i.e., low-poly, mesh, which can be easily textured. The mesh edges reproduce *sharp* details in the input (e.g., edges of the legs), yet still approximate smooth geometry (e.g., circular table-top). (b) State-of-the-art methods regress an indicator function, which needs to be iso-surfaced, resulting in over-tessellated meshes which only *approximate* sharp details with smooth surfaces.

surface representation; they are more controllable, easier to manipulate, and are more *compact*, attaining higher visual quality using fewer primitives; see Figure 1.

For visualization purposes, the generated voxels, point clouds, and implicits are typically converted into meshes in post-processing, e.g., via iso-surface extraction by Marching Cubes [26]. Few deep networks can generate polygonal meshes directly, and such methods are limited to genus-zero meshes [17, 27, 44], piece-wise genus-zero [12] meshes, meshes sharing the same connectivity [11, 41], or meshes with very low number of vertices [7]. Patch-based approaches can generate results which cover a 3D shape with planar polygons [46] or curved [16] mesh patches, but their visual quality is often tampered by visible seams, incoherent patch connections, and rough surface appearance. It is difficult to texture or manipulate such mesh outputs.

In this paper, we develop a generative neural network which outputs polygonal meshes *natively*. Specifically, parameters or weights that are learned by the network can predict multiple planes which fit the surfaces of a 3D shape, resulting in a *compact* and *watertight* polygonal mesh; see Figure 1. We name our network *BSP-Net*, since each facet is associated with a *binary space partitioning* (BSP), and the shape is composed by combining these partitions.

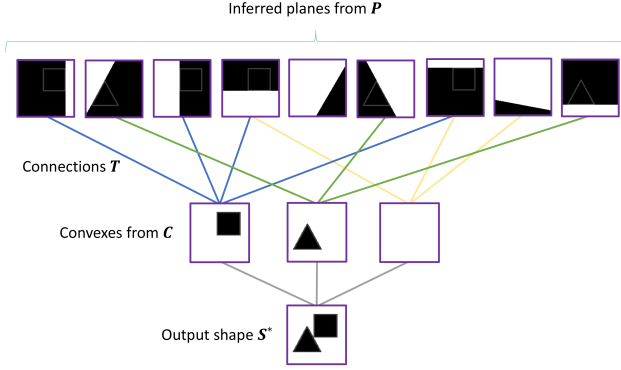


Figure 2: An illustration of “neural” BSP-tree.

BSP-Net learns an *implicit field*: given n point coordinates and a shape feature vector as input, the network outputs values indicating whether the points are inside or outside the shape. The construction of this implicit function is illustrated in Figure 2, and consists of three steps: ① a collection of plane equations implies a collection of p binary partitions of space; see Figure 2-top; ② an operator $\mathbf{T}_{p \times c}$ groups these partitions to create a collection of c convex shape primitives/parts; ③ finally, the part collection is merged to produce the implicit field of the output shape.

Figure 3 shows the network architecture of BSP-Net corresponding to these three steps: ① given the feature code, an MLP produces in layer L_0 a matrix $\mathbf{P}_{p \times 4}$ of canonical parameters that define the implicit equations of p planes: $ax + by + cz + d = 0$; these implicit functions are evaluated on a collection of n point coordinates $\mathbf{x}_{n \times 4}$ in layer L_1 ; ② the operator $\mathbf{T}_{p \times c}$ is a *binary* matrix that enforces a *selective* neuron feed from L_1 to the next network layer L_2 , forming convex parts; ③ finally, layer L_3 assembles the parts into a shape via either sum or min-pooling.

At inference time, we feed the input to the network to obtain components of the BSP-tree, i.e., leaf nodes (planes \mathbf{P}) and connections (binary weights \mathbf{T}). We then apply classic Constructive Solid Geometry (CSG) to extract the explicit polygonal surfaces of the shapes. The mesh is typically compact, formed by a subset of the p planes directly from the network, leading to a significant speed-up over the previous networks during inference, and without the need for expensive iso-surfacing – current inference time is about 0.5 seconds per generated mesh. Furthermore, meshes generated by the network are guaranteed to be watertight, possibly with *sharp* features, in contrast to smooth shapes produced by previous implicit decoders [5, 22, 28].

BSP-Net is trainable and characterized by *interpretable* network parameters defining the hyper-planes and their formation into the reconstructed surface. Importantly, the network training is *self-supervised* as no ground truth convex shape decompositions are needed. BSP-Net is trained to *reconstruct* all shapes from the training set using the *same* set

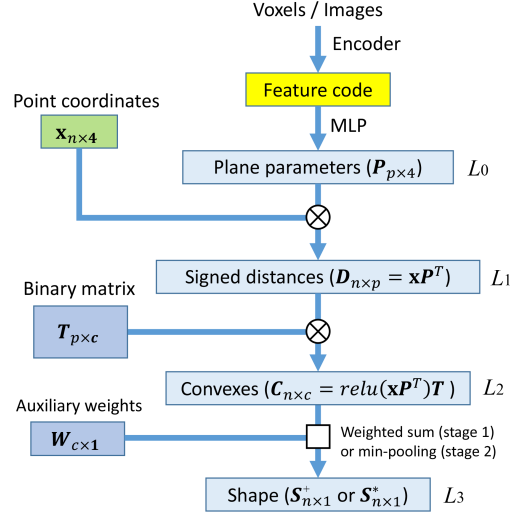


Figure 3: The network corresponding to Figure 2.

of convexes constructed in layer L_2 of the network. As a result, our network provides a *natural correspondence* between all the shapes at the level of the convexes. BSP-Net does not yet learn semantic parts. Grouping of the convexes into semantic parts can be obtained manually, or learned otherwise as semantic shape segmentation is a well-studied problem. Such a grouping only need to be done on each convex once to propagate the semantic understanding to all shapes containing the same semantic parts.

Contributions.

- BSP-Net is the first deep generative network which directly outputs compact and watertight polygonal meshes with arbitrary topology and structure variety.
- The learned BSP-tree allows us to infer both shape segmentation and part correspondence.
- By adjusting the encoder of our network, BSP-Net can also be adapted for shape auto-encoding and single-view 3D reconstruction (SVR).
- To the best of our knowledge, BSP-Net is among the first to achieve *structured* SVR, reconstructing a *segmented* 3D shape from a single unstructured object image.
- Last but not the least, our network is also the first which can reconstruct and recover sharp geometric features.

Through extensive experiments on shape auto-encoding, segmentation, part correspondence, and single-view reconstruction, we demonstrate state-of-the-art performances by BSP-Net. Comparisons are made to leading methods on shape decomposition and 3D reconstruction, using conventional distortion metrics, visual similarity, as well as a new metric assessing the capacity of a model in representing sharp features. In particular, we highlight the favorable fidelity-complexity trade-off exhibited by our network.

2. Related work

Large shape collections such as ShapeNet [2] and PartNet [30] have spurred the development of learning techniques for 3D data processing. In this section, we cover representative approaches based on the underlying shape representation learned, with a focus on generative models.

Grid models. Early approaches generalized 2D convolutions to 3D [6, 14, 24, 48, 49], and employed volumetric *grids* to represent shapes in terms of coarse occupancy functions, where a voxel evaluates to zero if it is outside and one otherwise. Unfortunately, these methods are typically limited to low resolutions of at most 64^3 due to the cubic growth in memory requirements. To generate finer results, differentiable marching cubes operations have been proposed [26], as well as hierarchical strategies [18, 36, 42, 45, 46] that alleviate the curse of dimensionality affecting dense volumetric grids. Another alternative is to use multi-view images [25, 40] and geometry images [38, 39], which allow standard 2D convolution, but such methods are only suitable on the *encoder* side of a network architecture, while we focus on decoders. Finally, recent methods that perform sparse convolutions [15] on voxel grids are similarly limited to encoders.

Surface models. As much of the semantics of 3D models is captured by their *surface*, the boundary between inside/outside space, a variety of methods have been proposed to represent shape surfaces in a differentiable way. Amongst these we find a category of techniques pioneered by PointNet [33] that express surfaces as point clouds [1, 9, 10, 33, 34, 53, 56], and techniques pioneered by AtlasNet [16] that adopt a 2D-to-3D mapping process [47, 39, 44, 53]. An interesting alternative is to consider mesh generation as the process of estimating vertices and their connectivity [7], but these methods do not guarantee watertight results, and hardly scale beyond a hundred vertices.

Implicit models. A very recent trend has been the modeling of shapes as a learnable indicator *function* [5, 22, 28], rather than a sampling of it, as in the case of voxel methods. The resulting networks treat reconstruction as a *classification* problem, and are universal approximators [20] whose reconstruction precision is proportional to the network complexity. However, at inference time, generating a 3D model still requires the execution of an expensive iso-surfacing operation whose performance scales cubically in the desired resolution. In contrast, our network directly outputs a *low-poly* approximation of the shape surface.

Shape decomposition. BSP-Net generates meshes using a part-based approach, hence techniques that learn shape decompositions are of particular relevance. There are methods that decompose shapes as oriented boxes [43, 31], axis aligned gaussians [13], super-quadratics [32], or a union of

indicator functions, in BAE-NET [4]. The architecture of our network draws inspiration from BAE-NET, which is designed to segment a shape by reconstructing its parts in different branches of the network. For each shape part, BAE-NET learns an implicit field by means of a binary classifier. In contrast, BSP-Net explicitly learns a tree structure built on plane subdivisions for bottom-up part assembly.

Another similar work is CvxNet [8], which decomposes shapes as a collection of convex primitives. However, BSP-Net differs from CvxNet in several significant ways: ① we target low-poly reconstruction with sharp features, while they target smooth reconstruction; ② their network always outputs K convexes, while the “right” number of primitives is learnt automatically in our method; ③ our optimization routine is completely different from theirs, as their compositional tree structure is *hard-coded*.

Structured models. There have been recent works on learning structured 3D models, in particular, linear [58] or hierarchical [23, 57, 29, 31] organization of part bounding boxes. While some methods learn part geometries separately [23, 29], others jointly embed/encode structure and geometry [51, 12]. What is common about all of these methods is that they are *supervised*, and were trained on shape collections with part segmentations and labels. In contrast, BSP-Net is unsupervised. On the other hand, our network is not designed to infer shape semantics; it is trained to learn convex decompositions. To the best of our knowledge, there is only one prior work, Im2Struct [31], which infers part structures from a single-view image. However, this work only produces a box arrangement; it does not reconstruct a structured *shape* like BSP-Net.

Binary and capsule networks. The discrete optimization for the tree structures in BSP-Net bears some resemblance to binary [21] and XNOR [35] neural networks. However, only *one* layer of BSP-Net employs binary weights, and our training method differs, as we use a continuous relaxation of the weights in early training. Further, as our network can be thought of as a simplified scene graph, it holds striking similarities to the principles of capsule networks [37], where low-level capsules (hyperplanes) are aggregated in higher (convexes) and higher (shapes) capsule representations. Nonetheless, while [37] addresses discriminative tasks (encoder), we focus on generative tasks (decoder).

3. Method

We seek a deep representation of geometry that is simultaneously trainable and interpretable. We achieve this task by devising a network architecture that provides a differentiable Binary Space Partitioning tree (BSP-tree) representation¹. This representation is easily *trainable* as it encodes

¹While typical BSP-trees are binary, we focus on n -ary trees, with the “B” in BSP referring to binary space partitioning, not the tree structure.

geometry via implicit functions, and *interpretable* since its outputs are a collection of convex polytopes. While we generally target 3D geometry, we employ 2D examples to explain the technique without loss of generality.

We achieve our goal via a network containing three main modules, which act on feature vectors extracted by an encoder corresponding to the type of input data (e.g. the features produced by ResNet for images or 3D CNN for voxels). In more details, a first layer that *extracts* hyperplanes conditional on the input data, a second layer that *groups* hyperplanes in the form of half-spaces to create parts (convexes), and a third layer *assembles* parts together to reconstruct the overall object; see Figure 3.

Layer 1: hyperplane extraction. Given a feature vector \mathbf{f} , we apply a multi-layer perceptron \mathcal{P} to obtain plane parameters $\mathbf{P}_{p \times 4}$, where p is the number of planes – i.e. $\mathbf{P} = \mathcal{P}_\omega(\mathbf{f})$. For any point $\mathbf{x} = (x, y, z, 1)$, the product $\mathbf{D} = \mathbf{x}\mathbf{P}^T$ is a vector of *signed* distances to each plane – the i th distance is negative if \mathbf{x} is *inside*, and positive if it is *outside*, the i th plane, with respect to the plane normal.

Layer 2: hyperplane grouping. To group hyperplanes into geometric primitives we employ a binary matrix $\mathbf{T}_{p \times c}$. Via a max-pooling operation we aggregate input planes to form a set of c *convex* primitives:

$$C_j^*(\mathbf{x}) = \max_i (D_i T_{ij}) \begin{cases} < 0 & \text{inside} \\ > 0 & \text{outside.} \end{cases} \quad (1)$$

Note that during training the gradients would flow through only one (max) of the planes. Hence, to ease training, we employ a version that replaces max with summation:

$$C_j^+(\mathbf{x}) = \sum_i \text{relu}(D_i) T_{ij} \begin{cases} = 0 & \text{inside} \\ > 0 & \text{outside.} \end{cases} \quad (2)$$

Layer 3: shape assembly. This layer groups convexes to create a possibly non-convex output shape via min-pooling:

$$S^*(\mathbf{x}) = \min_j (C_j^*(\mathbf{x})) \begin{cases} = 0 & \text{inside} \\ > 0 & \text{outside.} \end{cases} \quad (3)$$

Note that the use of C^+ in the expression above is *intentional*. We avoid using C^* due to the lack of a memory efficient implementation of the operator in TensorFlow 1.

Again, to facilitate learning, we distribute gradients to all convexes by resorting to a (weighted) summation:

$$S^+(\mathbf{x}) = \left[\sum_j W_j [1 - C_j^+(\mathbf{x})]_{[0,1]} \right]_{[0,1]} \begin{cases} = 1 & \approx \text{in} \\ [0, 1) & \approx \text{out,} \end{cases} \quad (4)$$

where $\mathbf{W}_{c \times 1}$ is a weight vector, and $[\cdot]_{[0,1]}$ performs clipping. During training we will enforce $\mathbf{W} \approx \mathbf{1}$. Note that

the inside/outside status here is only *approximate*. For example, when $\mathbf{W} = \mathbf{1}$, and all $C_j^+ = 0.5$, one is outside of all convexes, but inside their composition.

Two-stage training. Losses evaluated on (4) will be approximate, but have better gradient than (3). Hence, we develop a two-stage training scheme where: ① in the *continuous* phase, we try to keep all weights continuous and compute an approximate solution via $S^+(\mathbf{x})$ – this would generate an approximate result as can be observed in Figure 4 (b); ② in the next *discrete* phase, we quantize the weights and use a perfect union to generate accurate results by fine-tuning on $S^*(\mathbf{x})$ – this creates a much finer reconstruction as illustrated in Figure 4 (c,d).

Our two-stage training strategy is inspired by classical optimization, where smooth relaxation of integer problems is widely accepted, and mathematically principled.

3.1. Training Stage 1 – Continuous

We initialize \mathbf{T} and \mathbf{W} with random zero-mean Gaussian noise having $\sigma = 0.02$, and optimize the network via:

$$\arg \min_{\omega, \mathbf{T}, \mathbf{W}} \mathcal{L}_{\text{rec}}^+ + \mathcal{L}_{\mathbf{T}}^+ + \mathcal{L}_{\mathbf{W}}^+. \quad (5)$$

Given query points \mathbf{x} , our network is trained to match $S(\mathbf{x})$ to the ground truth indicator function, denoted by $F(\mathbf{x}|G)$, in a least-squares sense:

$$\mathcal{L}_{\text{rec}}^+ = \mathbb{E}_{\mathbf{x} \sim G} [(S^+(\mathbf{x}) - F(\mathbf{x}|G))^2], \quad (6)$$

where $\mathbf{x} \sim G$ indicates a sampling that is specific to the training shape G – including random samples in the unit box as well as samples near the boundary ∂G ; see [5]. An edge between plane i and convex j is represented by $\mathbf{T}_{ij} = 1$, and the entry is zero otherwise. We perform a continuous relaxation of a graph adjacency matrix \mathbf{T} , where we require its values to be bounded in the $[0, 1]$ range:

$$\mathcal{L}_{\mathbf{T}}^+ = \sum_{t \in \mathbf{T}} \max(-t, 0) + \sum_{t \in \mathbf{T}} \max(t - 1, 0). \quad (7)$$

Note that this is more effective than using a sigmoid activation, as its gradients do not vanish. Further, we would like \mathbf{W} to be close to $\mathbf{1}$ so that the merge operation is a sum:

$$\mathcal{L}_{\mathbf{W}}^+ = \sum_j |W_j - 1|. \quad (8)$$

However, we remind the reader that we initialize with $\mathbf{W} \approx \mathbf{0}$ to avoid vanishing gradients in early training.

3.2. Training Stage 2 – Discrete

In the second stage, we first quantize \mathbf{T} by picking a threshold $\lambda = 0.01$ and assign $t = (t > \lambda) ? 1 : 0$. Experimentally, we found the values learnt for \mathbf{T} to be small, which led

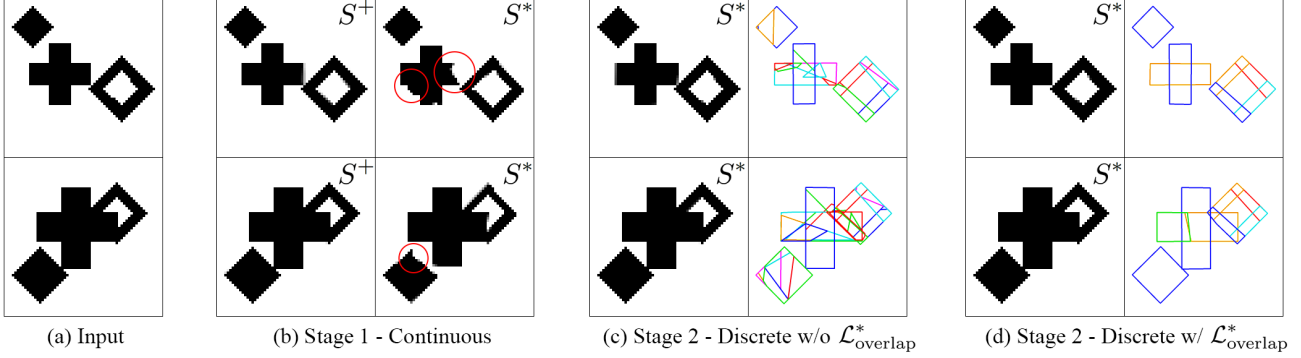


Figure 4: **Evaluation in 2D** – auto-encoder trained on the synthetic 2D dataset. We show auto-encoding results and highlight mistakes made in Stage 1 with red circles, which are resolved in Stage 2. We further show the effect of enabling the (optional) overlap loss. Notice that in the visualization we use different (possibly repeating) colors to indicate different convexes.

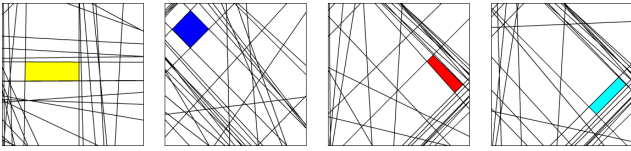


Figure 5: **Examples of L_2 output** – a few convexes from the first shape in Figure 4, and the planes to construct them. Note how many planes are unused.

to our choice of a small threshold value. With the quantized \mathbf{T} , we fine-tune the network by:

$$\arg \min_{\omega} \mathcal{L}_{\text{recon}}^* + \mathcal{L}_{\text{overlap}}^*, \quad (9)$$

where we ensure that the shape is well reconstructed via:

$$\mathcal{L}_{\text{recon}}^* = \mathbb{E}_{\mathbf{x} \sim G} [\mathbf{F}(\mathbf{x}|G) \cdot \max(S^*(\mathbf{x}), 0)] \quad (10)$$

$$+ \mathbb{E}_{\mathbf{x} \sim G} [(1 - \mathbf{F}(\mathbf{x}|G)) \cdot (1 - \min(S^*(\mathbf{x}), 1))]. \quad (11)$$

The above loss function pulls $S^*(\mathbf{x})$ towards 0 if \mathbf{x} should be inside the shape; it pushes $S^*(\mathbf{x})$ beyond 1 otherwise. Optionally, we can also discourage overlaps between the convexes. We first compute a mask M such that $M(\mathbf{x})=1$ if \mathbf{x} is contained in *more* than one convex, and then evaluate:

$$\mathcal{L}_{\text{overlap}}^* = -\mathbb{E}_{\mathbf{x} \sim G} [M(\mathbf{x})S^*(\mathbf{x})]. \quad (12)$$

3.3. Algorithmic and training details

In our 2D experiments, we use $p=256$ planes and $c=64$ convexes. We use a simple 2D convolutional encoder where each layer downsamples the image by half, and doubles the number of feature channels. We use the centers of all pixels as samples. In our 3D experiments, we use $p=4,096$ planes and $c=256$ convexes. The encoder for *voxels* is a 3D CNN encoder where each layer downsamples the grid by half, and doubles the number of feature channels. It takes a volume of size 64^3 as input. The encoder for *images* is ResNet-18 without pooling layers that receives images of size 128^2 as

input. All encoders produce feature codes $|\mathbf{f}|=256$. The dense network \mathcal{P}_{ω} has widths $\{512, 1024, 2048, 4p\}$ where the last layer outputs the plane parameters.

When training the auto-encoder for 3D shapes, we adopt the progressive training from [5], on points sampled from grids that are increasingly denser ($16^3, 32^3, 64^3$). Note that the hierarchical training is not necessary for convergence, but results in an $\approx 3\times$ speedup in convergence. In Stage 1, we train the network on 16^3 grids for 8 million iterations with batch size 36, then 32^3 for 8 million iterations with batch size 36, then 64^3 for 8 million iterations with batch size 12. In Stage 2, we train the network on 64^3 grids for 8 million iterations with batch size 12.

For single-view reconstruction, we also adopt the training scheme in [5], i.e., train an auto-encoder first, then only train the image encoder of the SVR model to predict *latent* codes instead of directly predicting the output shapes. We train the image encoder for 1,000 epochs with batch size 64. We run our experiments on a workstation with an Nvidia GeForce RTX 2080 Ti GPU. When training the auto-encoder (one model on the 13 ShapeNet categories), Stage 1 takes about ≈ 3 days and Stage 2 takes ≈ 2 days; training the image-encoder requires ≈ 1 day.

4. Results and evaluation

We study the behavior of BSP-Net on a synthetic 2D shape dataset (Section 4.1), and evaluate our auto-encoder (Section 4.2), as well as single view reconstruction (Section 4.3) compared to other state-of-the-art methods.

4.1. Auto-encoding 2D shapes

To illustrate how our network works, we created a synthetic 2D dataset. We place a diamond, a cross, and a hollow diamond with varying sizes over 64×64 images; see Figure 4(a). The order of the three shapes is *sorted* so that the diamond is always on the left and the hollow diamond is always on the right – this is to *mimic* the structure of shape

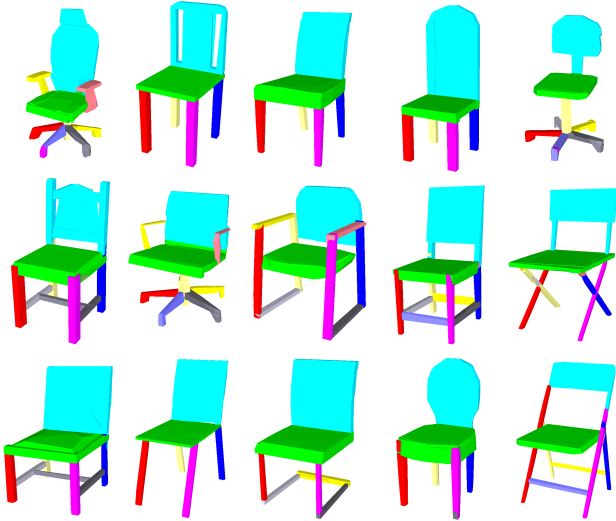


Figure 6: **Segmentation and correspondence** – semantics implied from autoencoding by BSP-Net. Colors shown here are the result of a *manual* grouping of learned convexes. The color assignment was performed on a few shapes: once a convex is colored in one shape, we can propagate the color to the other shapes by using the learnt convex id.

	CD	NC	LFD
VP [43]	2.259	0.683	6132.74
SQ [32]	1.656	0.719	5451.44
BAE [4]	1.592	0.777	4587.34
Ours	0.447	0.858	2019.26
Ours + $\mathcal{L}_{\text{overlap}}^*$	0.448	0.858	2030.35

Table 1: **Surface reconstruction** quality and comparison for 3D shape autoencoding. Best results are marked in bold.

	plane	car	chair	lamp	table	mean
VP [43]	37.6	41.9	64.7	62.2	62.1	56.9
SQ [32]	48.9	49.5	65.6	68.3	77.7	66.2
BAE [4]	40.6	46.9	72.3	41.6	68.2	59.8
Ours	74.2	69.5	80.9	52.3	90.3	79.3
Ours + $\mathcal{L}_{\text{overlap}}^*$	74.5	69.7	82.1	53.4	90.3	79.8
BAE* [4]	75.4	73.5	85.2	73.9	86.4	81.8

Table 2: **Segmentation**: comparison in per-label IoU.

datasets such as ShapeNet [2]. After training Stage 1, our network has already achieved a good approximate S^+ reconstruction, however, by inspecting S^* , the output of our inference, we can see there are several imperfections. After the fine-tuning in Stage 2, our network achieves near perfect reconstructions. Finally, the use of overlap losses significantly improves the compactness of representation, reducing the number of convexes per part; see Figure 4(d).

Figure 5 visualizes the planes used to construct the individual convexes – we visualize planes i in convex j so that $T_{ij}=1$ and $P_{i1}^2 + P_{i2}^2 + P_{i3}^2 > \epsilon$ for a small threshold ϵ (to ignore planes with near-zero gradients). Note how BSP-Net

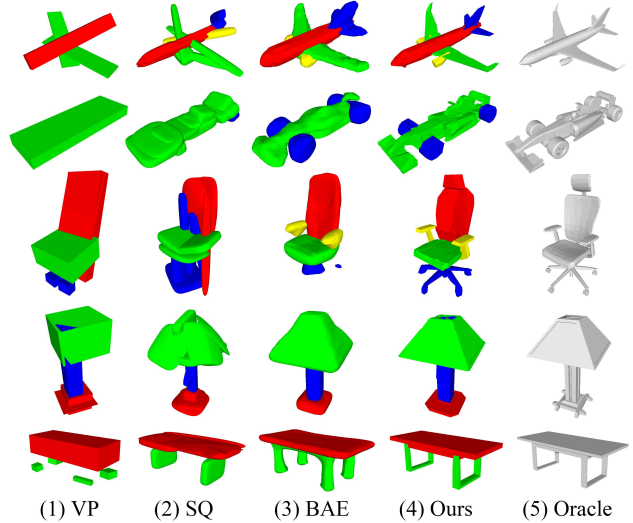


Figure 7: **Segmentation and reconstruction / Qualitative**.

creates a natural *semantic* correspondence across inferred convexes. For example, the hollow diamond in Figure 4(d) is always made of the same four convexes in the same relative positions – this is mainly due to the static structure in \mathcal{T} : different shapes need to *share* the same set of convexes and their associated hyper-planes.

4.2. Auto-encoding 3D shapes

For 3D shape autoencoding, we compare BSP-Net to a few other shape decomposition networks: Volumetric Primitives (VP) [43], Super Quadrics (SQ) [32], and Branched Auto Encoders (BAE) [4]. Note that for the segmentation task, we also evaluate on BAE*, the version of BAE that uses the values of the predicted implicit function, and not just the classification boundaries – please note that the surface *reconstructed* by BAE and BAE* are *identical*.

Since all these methods target *shape decomposition* tasks, we train *single* class networks, and evaluate segmentation as well as reconstruction performance. We use the ShapeNet (Part) Dataset [54], and focus on five classes: airplane, car, chair, lamp and table. For the car class, since none of the networks separates surfaces (as we perform *vol-umetric* modeling), we reduce the parts from (wheel, body, hood, roof) \rightarrow (wheel, body); and analogously for lamps (base, pole, lampshade, canopy) \rightarrow (base, pole, lampshade) and tables (top, leg, support) \rightarrow (top, leg).

As quantitative metrics for reconstruction tasks, we report symmetric Chamfer Distance (**CD**, scaled by $\times 1000$) and Normal Consistency (**NC**) computed on $4k$ surface sampled points. We also report the Light Field Distance (**LFD**) [3] – the best-known visual similarity metric from computer graphics. For segmentation tasks, we report the typical mean per-label Intersection Over Union (IOU).

Segmentation. Table 2 shows the per category segmentation results. As we have ground truth part labels for the point clouds in the dataset, after training each network, we obtain the part label for each primitive/convex by *voting*: for each point we identify the nearest primitive to it, and then the point will cast a vote for that primitive on the corresponding part label. Afterwards, for each primitive, we assign to it the part label that has the highest number of votes. We use 20% of the dataset for assigning part labels, and we use all the shapes for testing. At test time, for each point in the *point cloud*, we find its nearest primitive, and assign the part label of the primitive to the point. In the comparison to BAE, we employ their one-shot training scheme [4, Sec.3.1]. Note that BAE-NET* is specialized to the segmentation task, while our work mostly targets part-based reconstruction; as such, the IoU performance in Table 2 is an *upper bound* of segmentation performance.

Figure 6 shows *semantic* segmentation and part correspondence implied by BSP-Net autoencoding, showing how individual parts (left/right arm/leg, etc.) are matched. In our method, all shapes are corresponded at the primitive (convexes) level. To reveal shape semantics, we *manually* group convexes belonging to the same semantic part and assign them the same color. Note that the color assignment is done on each convex once, and propagated to all the shapes.

Reconstruction comparison. BSP-Net achieves significantly better reconstruction quality, while maintaining high segmentation accuracy; see Table 1 and Figure 7, where we color each *primitive* based on its inferred part label. BAE-NET was designed for segmentation, thus produces poor-quality part-based 3D reconstructions. Note how BSP-Net is able to represent complex parts such as legs of swivel chairs in Figure 7, while none of the other methods can.

4.3. Single view reconstruction (SVR)

We compare our method with AtlasNet [16], IM-NET [5] and OccNet [28] on the task of single view reconstruction. We report quantitative results in Table 3 and Table 4, and qualitative results in Figure 8. We use the 13 categories in ShapeNet [2] that have more than 1,000 shapes each, and the rendered views from 3D-R2N2 [6]. We train one model on all categories, using 80% of the shapes for training and 20% for testing, in a similar fashion to AtlasNet [16]. For other methods, we download the *pre-trained* models released by the authors. Since the pre-trained OccNet [28] model has a different train-test split than others, we evaluate it on the intersection of the test splits.

Edge Chamfer Distance (ECD). To measure the capacity of a model to represent *sharp* features, we introduce a new metric. We first compute an “edge sampling” of the surface by generating $16k$ points $\mathbf{S}=\{\mathbf{s}_i\}$ uniformly distributed on the surface of a model, and then compute sharpness as:

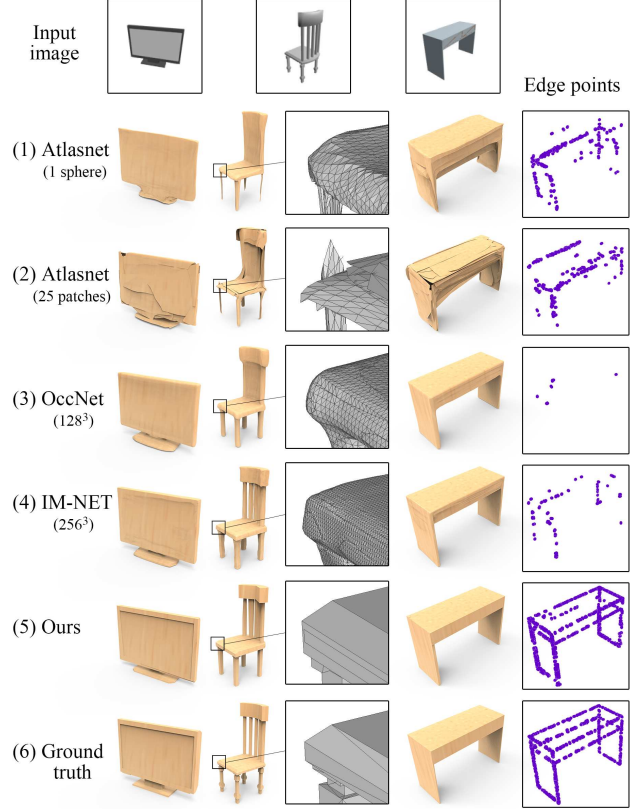


Figure 8: **Single-view 3D reconstruction** – comparison to AtlasNet [16], IM-NET [5], and OccNet [28]. Middle column shows mesh tessellations of the reconstruction; last column shows the edge sampling used in the ECD metric.

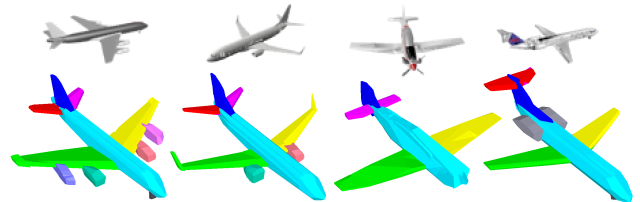


Figure 9: **Structured SVR** by BSP-Net reconstructs each shape with corresponding convexes. Convexes belonging to the same semantic parts are manually grouped and assigned the same color, resulting in semantic part correspondence.

$\sigma(\mathbf{s}_i) = \min_{j \in \mathcal{N}_\varepsilon(\mathbf{s}_i)} |\mathbf{n}_i \cdot \mathbf{n}_j|$, where $\mathcal{N}_\varepsilon(\mathbf{s})$ extracts the indices of the samples in \mathbf{S} within distance ε from \mathbf{s} , and \mathbf{n} is the surface normal of a sample. We set $\varepsilon=0.01$, and generate our edge sampling by retaining points such that $\sigma(\mathbf{s}_i) < 0.1$; see Figure 8. Given two shapes, the ECD between them is nothing but the Chamfer Distance between the corresponding edge samplings.

Analysis. Our method achieves *comparable* performance to the state-of-the-art in terms of Chamfer Distance. As for visual quality, our method also *outperforms* most other

	Chamfer Distance (CD)					Edge Chamfer Distance (ECD)					Light Field Distance (LFD)				
	Atlas0	Atlas25	OccNet ₃₂	IM-NET ₃₂	Ours	Atlas0	Atlas25	OccNet ₃₂	IM-NET ₃₂	Ours	Atlas0	Atlas25	OccNet ₃₂	IM-NET ₃₂	Ours
airplane	0.587	0.440	1.534	2.211	0.759	0.396	0.575	1.494	0.815	0.487	5129.36	4680.37	7760.42	7581.13	4496.91
bench	1.086	0.888	3.220	1.933	1.226	0.658	0.857	2.131	1.400	0.475	4387.28	4220.10	4922.89	4281.18	3380.46
cabinet	1.231	1.173	1.099	1.902	1.188	3.676	2.821	10.804	9.521	0.435	1369.90	1558.45	1187.08	1347.97	989.12
car	0.799	0.688	0.870	1.390	0.841	1.385	1.279	8.428	6.085	0.702	1870.42	1754.87	1790.00	1932.78	1694.81
chair	1.629	1.258	1.484	1.783	1.340	1.440	1.951	4.262	3.545	0.872	3993.94	3625.23	3354.00	3473.62	2961.20
display	1.516	1.285	2.171	2.370	1.856	2.267	2.911	6.059	5.509	0.697	2940.36	3004.44	2565.07	3232.06	2533.86
lamp	3.858	3.248	12.528	6.387	3.480	2.458	2.690	8.510	4.308	2.144	7566.25	7162.20	8038.98	6958.52	6726.92
speaker	2.328	1.957	2.662	3.120	2.616	9.199	5.324	11.271	9.889	1.075	2054.18	2075.69	2393.50	1955.40	1748.26
rifle	1.001	0.715	2.015	2.052	0.888	0.288	0.318	1.463	1.882	0.231	6162.03	6124.89	6615.20	6070.86	4741.70
couch	1.471	1.233	1.246	2.344	1.645	2.253	3.817	10.179	8.531	0.869	2387.09	2343.11	1956.26	2184.28	1880.21
table	1.996	1.376	3.734	2.778	1.643	1.122	1.716	3.900	3.097	0.515	3598.59	3286.05	3371.20	3347.12	2627.82
phone	1.048	0.975	1.183	2.268	1.383	10.459	11.585	16.021	14.684	1.477	1817.61	1816.22	1995.98	1964.46	1555.47
vessel	1.179	0.966	1.691	2.385	1.585	0.782	0.889	12.375	3.253	0.588	4551.17	4430.04	5066.99	4494.14	3931.73
mean	1.487	1.170	2.538	2.361	1.432	1.866	2.069	6.245	4.617	0.743	3644.91	3436.14	3795.23	3700.22	2939.15

Table 3: **Single view reconstruction** – comparison to the state of the art. *Atlas25* denotes AtlasNet with 25 square patches, while *Atlas0* uses a single spherical patch. Subscripts to OccNet and IM-NET show sampling resolution. For fair comparisons, we use resolution 32^3 so that OccNet and IM-NET output meshes with comparable number of vertices and faces.

	CD	ECD	LFD	#V	#F
Atlas0	1.487	1.866	3644.91	7446	14888
Atlas25	1.170	2.069	3436.14	2500	4050
OccNet ₃₂	2.538	6.245	3795.23	1511	3017
OccNet ₆₄	1.950	6.654	3254.55	6756	13508
OccNet ₁₂₈	1.945	6.766	3224.33	27270	54538
IM-NET ₃₂	2.361	4.617	3700.22	1204	2404
IM-NET ₆₄	1.467	4.426	2940.56	5007	10009
IM-NET ₁₂₈	1.387	1.971	2810.47	20504	41005
IM-NET ₂₅₆	1.371	2.273	2804.77	82965	165929
Ours	1.432	0.743	2939.15	1191	1913

Table 4: **Low-poly analysis** – the dataset-averaged metrics in single view reconstruction. We highlight the number of vertices #V and triangles #F in the predicted meshes.

methods, which is reflected by the superior results in terms of Light Field Distance. Similarly to Figure 6, we manually color each convex to show part correspondences in Figure 9. We visualize the triangulations of the output meshes in Figure 8: our method outputs meshes with a smaller number of polygons than state-of-the-art methods. Note that these methods cannot generate low-poly meshes, and their vertices are always distributed quasi-uniformly.

Finally, note that our method is the *only* one amongst those tested capable of representing sharp edges – this can be observed quantitatively in terms of Edge Chamfer Distance, where BSP-Net performs much better. Note that AtlasNet could also generate edges in theory, but the shape is not watertight and the edges are irregular, as it can be seen in the zoom-ins of Figure 8. We also analyze these metrics aggregated on the entire testing set in Table 4. In this final analysis, we also include OccNet₁₂₈ and IM-NET₂₅₆, which are the *original* resolutions used by the authors. Note the average number of *polygons* inferred by our method is 655 (recall #polygons \leq #triangles in polygonal meshes).

5. Conclusion, limitation, and future work

We introduce BSP-Net, an unsupervised method which can generate compact and structured polygonal meshes in

the form of convex decomposition. Our network learns a BSP-tree built on the same set of planes, and in turn, the same set of convexes, to minimize a reconstruction loss for the training shapes. These planes and convexes are defined by weights learned by the network. Compared to state-of-the-art methods, meshes generated by BSP-Net exhibit superior visual quality, in particular, sharp geometric details, when *comparable* number of primitives are employed.

The main limitation of BSP-Net is that it can only decompose a shape as a *union* of convexes. Concave shapes, e.g., a teacup or ring, have to be decomposed into many small convex pieces, which is unnatural and leads to wasting of a considerable amount of representation budget (planes and convexes). A better way to represent such shapes is to do a *difference* operation rather than union. How to generalize BSP-Net to express a variety of CSG operations is an interesting direction for future work.

Current training times for BSP-Net are quite significant: 6 days for 4,096 planes and 256 convexes for the SVR task trained across all categories; inference is fast however. While most shapes only need a small number of planes to represent, we cannot reduce the total number of planes as they are needed to well represent a large *set* of shapes. It would be ideal if the network can adapt the primitive count based on the complexity of the input shapes; this may call for an architectural change to the network.

While its applicability to RGBD data could leverage the auto-decoder ideas explored by [22], the generalization of our method beyond curated datasets [2], and the ability to train from only RGB images are of critical importance.

Acknowledgements

We thank the anonymous reviewers for their valuable comments. This work was supported in part by an NSERC grant (No. 611370), a Google Faculty Research Award, and Google Cloud Platform research credits.

References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J Guibas. Learning representations and generative models for 3d point clouds. In *ICML*, 2018. 1, 3
- [2] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An information-rich 3D model repository. *arXiv preprint arXiv:1512.03012*, 2015. 3, 6, 7, 8
- [3] Ding-Yun Chen, Xiao-Pei Tian, Yu-Te Shen, and Ming Ouhyoung. On visual similarity based 3d model retrieval. In *Computer graphics forum*, 2003. 6
- [4] Zhiqin Chen, Kangxue Yin, Matthew Fisher, Siddhartha Chaudhuri, and Hao Zhang. BAE-NET: Branched autoencoder for shape co-segmentation. *ICCV*, 2019. 3, 6, 7
- [5] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. *CVPR*, 2019. 1, 2, 3, 4, 5, 7
- [6] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *ECCV*, 2016. 3, 7
- [7] Angela Dai and Matthias Nießner. Scan2mesh: From unstructured range scans to 3d meshes. In *CVPR*, 2019. 1, 3
- [8] Boyang Deng, Kyle Genova, Soroosh Yazdani, Sofien Bouaziz, Geoffrey Hinton, and Andrea Tagliasacchi. Cvxnets: Learnable convex decomposition. In *CVPR*, 2020. 3
- [9] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *CVPR*, 2017. 3
- [10] Matheus Gadelha, Rui Wang, and Subhransu Maji. Multiresolution tree networks for 3d point cloud processing. In *ECCV*, 2018. 3
- [11] Lin Gao, Yu-Kun Lai, Jie Yang, Ling-Xiao Zhang, Leif Kobbelt, and Shihong Xia. Sparse data driven mesh deformation. *IEEE transactions on visualization and computer graphics*, 2019. 1
- [12] Lin Gao, Jie Yang, Tong Wu, Hongbo Fu, Yu-Kun Lai, and Hao Zhang. SDM-NET: Deep generative network for structured deformable mesh. *ACM Trans. on Graphics (TOG)*, 2019. 1, 3
- [13] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T. Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *ICCV*, 2019. 1, 3
- [14] Rohit Girdhar, David F Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and generative vector representation for objects. In *ECCV*, 2016. 1, 3
- [15] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *CVPR*, 2018. 3
- [16] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan Russell, and Mathieu Aubry. Atlasnet: A papier-mâché approach to learning 3d surface generation. In *CVPR*, 2018. 1, 3, 7
- [17] Heli Ben Hamu, Haggai Maron, Itay Kezurer, Gal Avineri, and Yaron Lipman. Multi-chart generative surface modeling. *ACM Trans. on Graphics (TOG)*, 2018. 1
- [18] Christian Häne, Shubham Tulsiani, and Jitendra Malik. Hierarchical surface prediction for 3d object reconstruction. In *3DV*, 2017. 1, 3
- [19] Rana Hanocka, Amir Hertz, Noa Fish, Raja Giryes, Shachar Fleishman, and Daniel Cohen-Or. MeshCNN: A network with an edge. *ACM Trans. on Graphics (TOG)*, 2019. 1
- [20] Kurt Hornik. Approximation capabilities of multilayer feed-forward networks. *Neural Networks*, 1991. 3
- [21] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. In *NeurIPS*, 2016. 3
- [22] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019. 1, 2, 3, 8
- [23] Jun Li, Kai Xu, Siddhartha Chaudhuri, Ersin Yumer, Hao Zhang, and Leonidas Guibas. Grass: Generative recursive autoencoders for shape structures. *ACM Trans. on Graphics (TOG)*, 2017. 3
- [24] Yiyi Liao, Simon Donné, and Andreas Geiger. Deep marching cubes: Learning explicit surface representations. In *CVPR*, 2018. 3
- [25] Chen-Hsuan Lin, Chen Kong, and Simon Lucey. Learning efficient point cloud generation for dense 3d object reconstruction. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2018. 3
- [26] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. *SIGGRAPH*, 1987. 1, 3
- [27] Haggai Maron, Meirav Galun, Noam Aigerman, Miri Trope, Nadav Dym, Ersin Yumer, Vladimir G. Kim, and Yaron Lipman. Convolutional neural networks on surfaces via seamless toric covers. *ACM Trans. on Graphics (TOG)*, 2017. 1
- [28] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *CVPR*, 2019. 1, 2, 3, 7
- [29] Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas J Guibas. StructureNet: Hierarchical graph networks for 3d shape generation. *SIGGRAPH Asia*, 2019. 3
- [30] Kaichun Mo, Shilin Zhu, Angel X. Chang, Li Yi, Subarna Tripathi, Leonidas J. Guibas, and Hao Su. PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding. In *CVPR*, 2019. 3
- [31] Chengjie Niu, Jun Li, and Kai Xu. Im2struct: Recovering 3d shape structure from a single RGB image. In *CVPR*, 2018. 3
- [32] Despoina Paschalidou, Ali Osman Ulusoy, and Andreas Geiger. Superquadrics revisited: Learning 3d shape parsing beyond cuboids. In *CVPR*, 2019. 3, 6
- [33] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *CVPR*, 2017. 1, 3

- [34] Charles R. Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *NeurIPS*, 2017. 1, 3
- [35] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *ECCV*, 2016. 3
- [36] Gernot Riegler, Ali Osman Ulusoy, Horst Bischof, and Andreas Geiger. Octnetfusion: Learning depth fusion from data. In *3DV*, 2017. 3
- [37] Sara Sabour, Nicholas Frosst, and Geoffrey E. Hinton. Dynamic routing between capsules. In *NeurIPS*, 2017. 3
- [38] Ayan Sinha, Jing Bai, and Karthik Ramani. Deep learning 3d shape surfaces using geometry images. In *ECCV*, 2016. 3
- [39] Ayan Sinha, Asim Unmesh, Qixing Huang, and Karthik Ramani. Surfnet: Generating 3d shape surfaces using deep residual networks. In *CVPR*, 2017. 3
- [40] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *ICCV*, 2015. 3
- [41] Qingyang Tan, Lin Gao, Yu-Kun Lai, and Shihong Xia. Variational autoencoders for deforming 3D mesh models. In *CVPR*, 2018. 1
- [42] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *ICCV*, 2017. 3
- [43] Shubham Tulsiani, Hao Su, Leonidas J Guibas, Alexei A Efros, and Jitendra Malik. Learning shape abstractions by assembling volumetric primitives. In *CVPR*, 2017. 3, 6
- [44] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *ECCV*, 2018. 1, 3
- [45] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-CNN: Octree-based Convolutional Neural Networks for 3D Shape Analysis. *SIGGRAPH*, 2017. 3
- [46] Peng-Shuai Wang, Chun-Yu Sun, Yang Liu, and Xin Tong. Adaptive O-CNN: A Patch-based Deep Representation of 3D Shapes. *SIGGRAPH Asia*, 2018. 1, 3
- [47] Francis Williams, Teseo Schneider, Claudio Silva, Denis Zorin, Joan Bruna, and Daniele Panozzo. Deep geometric prior for surface reconstruction. In *CVPR*, 2019. 3
- [48] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *NeurIPS*, 2016. 1, 3
- [49] Jiajun Wu, Chengkai Zhang, Xiuming Zhang, Zhoutong Zhang, William T. Freeman, and Joshua B. Tenenbaum. Learning shape priors for single-view 3d completion and reconstruction. In *ECCV*, 2018. 3
- [50] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *CVPR*, 2015. 1
- [51] Zhijie Wu, Xiang Wang, Di Lin, Dani Lischinski, Daniel Cohen-Or, and Hui Huang. SAGNet: Structure-aware generative network for 3D-shape modeling. *SIGGRAPH*, 2019. 3
- [52] Qiangeng Xu, Weiyue Wang, Duygu Ceylan, Radomír Mech, and Ulrich Neumann. DISN: deep implicit surface network for high-quality single-view 3d reconstruction. *NeurIPS*, 2019. 1
- [53] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *CVPR*, 2018. 3
- [54] Li Yi, Vladimir G Kim, Duygu Ceylan, I Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, Leonidas Guibas, et al. A scalable active framework for region annotation in 3D shape collections. 35(6), 2016. 6
- [55] Kangxue Yin, Zhiqin Chen, Hui Huang, Daniel Cohen-Or, and Hao Zhang. LOGAN: Unpaired shape transform in latent overcomplete space. *ACM Trans. on Graphics (TOG)*, 2019. 1
- [56] Kangxue Yin, Hui Huang, Daniel Cohen-Or, and Hao Zhang. P2P-NET: Bidirectional point displacement net for shape transform. *ACM Trans. on Graphics (TOG)*, 2018. 1, 3
- [57] Chenyang Zhu, Kai Xu, Siddhartha Chaudhuri, Renjiao Yi, and Hao Zhang. SCORES: Shape composition with recursive substructure priors. *ACM Trans. on Graphics (TOG)*, 2018. 3
- [58] Chuhan Zou, Ersin Yumer, Jimei Yang, Duygu Ceylan, and Derek Hoiem. 3D-PRNN: Generating shape primitives with recurrent neural networks. In *ICCV*, 2017. 3