

# Parameterization Invariant Learning on 3D Textured Meshes via Cross-Atlas Convolution

Anonymous CVPR submission

Paper ID 2492

## Abstract

*We present a convolutional network architecture for direct learning on textured meshes via texture maps, which encode the parameterization from 3D to 2D domain. The texture map can include not only RGB values but also rasterized geometric features if necessary. Since the parameterization of texture map could be unpredictable, we introduce the novel cross-atlas convolution to recover the original mesh geodesic neighborhood, so as to achieve the invariance property to arbitrary parameterization. The proposed module is integrated to classification and segmentation architectures, which take as input the texture map of a mesh, and output the predicted result. Our method not only shows competitive performances on public classification and segmentation benchmarks, but also unlocks a new practice to directly consume textured meshes.*

## 1. Introduction

The 3D mesh is one of the most popular representation for 3D shape, which consists of an array of vertices and an array of face indices indicating the surface geometry. It could be augmented with the texture coordinate array and the texture map to render the color appearance of the mesh.

Learning on textured meshes is challenging: on the geometry side, the arrays of vertices and faces are permutable. The mesh geometry by its design, could be very adaptive, meaning a similar shape can be meshed in very different patterns, with irregular vertex densities and inconsistent triangle qualities. On the texture side, the parameterization of texture map could be arbitrary and still renders the same appearance, as long as the texture coordinate of triangle is in accordance with the texture map.

Existing methods that can indirectly or directly operate on textured meshes include 1) multi-view projection [32, 11, 10], which renders the RGBD images from all around perspectives of the mesh. The multi-view images can be learned via individual 2D convolutional neural networks (CNNs) and the global feature is aggregated by view

pooling. However, this is only feasible for small objects [38] where occlusion is not significant. When it comes to larger scenes where occlusion and view selection are non-trivial, their methods would degrade. 2) Volumetric grids can be obtained from meshes, where voxels store the color value and then are learned by 3D CNNs [38, 5, 21]. However, as also mentioned in [6], the volumetric representation is memory-consuming and loses rich valuable image details, thus hindering the performance. 3) Point-based learning is a recent breakthrough [24, 26]. The colored point cloud can be easily sampled from the mesh. As the point cloud learning employs multi-layer perceptrons rather than convolution, it takes significantly more parameters and thus the maximum number of points can be learned is far behind image pixels under the same environment. 3) Geometric deep learning [3] techniques can directly process on meshes via spectral analyses [4, 8, 16, 13, 39] or geodesic convolutions [20, 2, 22]. While their methods focus more on learning the geometric shape for the dynamic correspondence task [1]. It's unclear how they can combine the texture or image for semantic learning.

In fact, the textured mesh is a self-contained combination between 2D texture and 3D shape. Some recent works are aware of this importance and perform joint learning on 2D images and 3D geometries [25, 6, 31], achieving improved performances. However, processing on images, rather than textures, is redundant since the same 3D area is seen and processed multiple times in images. Besides, such bundle of meshes, images and coresponding camera poses is difficult to acquire. With the popularity of 3D sensing techniques, it is more likely to obtain the standone textured mesh model.

To this end, we propose to perform direct learning on the textured mesh via its texture map. The texture map can include not only the color information as it usually does, but also arbitrary geometric features as long as they are rasterized in the map. Each pixel in texture map becomes a generic feature vector, encoding both color and geometric information. More importantly, the texture map is already in 2D domain, enjoying numerous benefits: 1) it can be learned via the standard CNNs, which can leverage the efficient de-

signs from rich previous researches. 2) The texture map rasterization is analogous to sampling, and thus invariant to the irregular meshing. 3) The hierarchy of a 2D map is simply the image pyramid, making multi-scale learning (*e.g.*, for global feature extraction) easily achievable. 4) The texture map inherits the geodesic neighborhood of meshes. *i.e.*, neighbor pixels on texture map indicate their corresponding 3D points must be neighbors on mesh surface, whereas other (volumetric or point-based) representations discard this neighborhood information.

The practice of learning 3D meshes via 2D domain is seen in previous methods [30, 19], but they have two crucial problems hindering their performances. The first one is distortion: unlike the generic texture map parameterization that segments the mesh to multiple atlases and pack them tightly in the texture map, their parameterization conducts only one cut on the mesh and unfolds it to a complete 2D squared map. This inevitably introduces distortions – when the mesh is far from a sphere, the distortion could be unpredictably large [30]. The second problem is their networks are not invariant to parameterization, which is further determined by the cut on the mesh surface. The author suggests to try multiple cuts on testing stage [19], and select the most responsive one, which makes it troublesome.

In this paper, we address above two problems. First, we do not unfold the 3D surface onto a complete 2D map. Instead, we use the generic UV parameterization which segments the mesh to multiple charts, projects them to 2D atlases and packs them as tight as possible inside the texture map. By doing so, each atlas finds its best projection direction to minimize the distortion. Second, the positions of atlases in texture map are unpredictable, depending on how the packing algorithm computes. Each atlas is isolated, meaning the neighborhood information is taken apart when crossing the texture seams on mesh surface. To tackle these issues, we introduce the *cross-atlas convolution*. When the filter convolves across the boundary of an atlas, the pixel located outside of the atlas is redirected to the correct position, which corresponds to the actual neighbor point on mesh surface. This is done with a precomputed offset map. We have integrated the cross-atlas convolution into the classification and segmentation network architects, and verified their effectiveness on public benchmarks. Overall, our method enjoys several benefits:

1. We unlocks the direct learning on textured meshes.
2. Our method addresses the distortion and the variance of parameterization problems in previous related methods [30, 19]. We also impose no restriction to the input mesh whereas they requires genus-0 meshes.
3. Our designed module is flexible and introduce no extra parameters: it can be trained on natural images, and applies inference on texture maps using our cross-atlas modules (see also Section 5.3).

## 2. Related Works

Deep learning on non-Euclidean geometric 3D data is an active and ongoing research topic. The mesh is one of the most commonly used representations in 3D vision and graphics, yet the irregularity of mesh makes it challenging to learn. We survey existing approaches in three categories.

**Converting to regular structures** The most straightforward approach is to convert the irregular mesh to regular data structure suitable for CNN processing. This can be done by projecting the mesh to multi-view 2D images, and then applying 2D CNN and view pooling to aggregate the global feature [32, 11, 10]. These methods show great performance on small object classification such as ModelNet [38] where viewing angles are simply from all around and the occlusion is not a problem, but their performance degrades when it comes to self-occluded objects or larger scenes where view selection is non-trivial. Another branch of methods is to convert the mesh to volumetric domains, and extract deep features from volumes by 3D convolutions [38, 5, 21]. The voxelization may introduce discretization errors and is highly memory-consuming. Using Octrees [27, 35] can abbreviate the resolution problem to some extent. Recently, some approaches combine both 2D views and 3D volumes and achieve even better results [25, 6].

**Point cloud approaches** Point cloud can be easily obtained by sampling on meshes. The irregular point cloud data can be learned from PointNet [24], and its extended hierarchical version [26]. They use combinations of multi-layer perceptions and pooling operations to achieve permutation invariance on the point set. Their insights also inspire several following works on point cloud learning in terms of improving the scalability and enhancing the local information of point cloud structure [17, 31, 14, 36]. In these methods, the neighborhood of a point is found by the radius-search or K nearest neighbors, which does not keep the original geodesic neighborhood on meshes. This could be a potential problem when the geodesic distance and Euclidean distance vary in the mesh model.

**Geometric deep learning on meshes** The mesh can be learned by geometric deep learning techniques [3] for the non-rigid shape correspondence task. 1) If the mesh is seen as a graph, several works have proposed to apply the spectral analysis on the eigen decomposition of the Laplacian of mesh graph [4, 8] to establish dense correspondences between deformable shapes. A general limitation is the cross-domain generalization issue, which is later addressed by spectral transformers [39] to some extent. Dirac operator is an alternative to Laplacian operator which yields better stability in some scenarios [16]. 2) If seen as a manifold

surface, the mesh can be learned by geodesic convolutions [20, 2, 22] on local patches, enjoying better generalization across domains than spectral methods. These techniques show great performance on non-rigid shape correspondence task but the receptive field of a polar filter is very small and thus it is unclear how to extract multi-scale features. Besides of using polar filter, a recent work [23] proposes to apply standard convolution on tangential projections of local patches and construct the hierarchy via mesh simplifications. 3) The third class of techniques applies global parameterization to the mesh and flatten it to 2D images [30, 19]. Our work belongs to this class, while the other two works are the GeometryImage [30] and the “flat-torus” method [19]. These methods enjoy common benefits derived from CNNs, but both of them unfold the mesh to a complete squared map, which induces considerable distortions. Besides, their methods require genus-0 input, otherwise they would crudely fill all topological holes. Regarding the parameterization, the one in GeometryImage [30] is not seamless, while the network in “flat-torus” method [19] is not invariant to the parameterization, depending on the cut of three chosen points on the mesh.

### 3. Mesh learning via Texture maps

This section illustrates the detailed procedure of mesh learning in the texture map space. In Section 3.1, we describe how the input texture map is generated from a pure mesh or a textured mesh. In Section 3.2, we introduce cross-atlas modules to recover the connectivities of separated atlases. In Section 3.3, we present the network architecture for classification and semantic segmentation tasks.

#### 3.1. Generating Texture Maps

In the preprocessing step, the input is a triangular mesh  $\mathcal{M} = \{V, T\}$ , where  $V = \{v_i\}$  and  $T = \{t_i\}$  correspond to the vertices and triangles respectively. If it's a polygon mesh we simply triangulate the faces. The mesh can be with or without texture coordinates and texture maps. The output includes a texture map ( $H \times W \times C$ ) and an offset map ( $H \times W \times 2k^2$ ) for the network input.

If the mesh is without textures (e.g., the CAD mesh in ModelNet [38], Fig. 2(a)), we create the UV parameterization for this mesh, and then rasterize the geometric features to a  $H \times W \times C$  texture map. The concrete algorithm is illustrated in Algorithm 1: we first find a minimum set of dominant projection directions  $\mathbf{P} = \{P_i \in \mathbb{R}^3\}$  such that the angle between each triangle normal  $\text{normal}(t_i)$  and its best projection vector  $P_{t_i}$  is less than a threshold, i.e.,  $\angle(\text{normal}(t_i), P_{t_i}) < \tau_{\text{angle}}$ . When the angle  $\tau_{\text{angle}}$  approaches  $0^\circ$ , the projection is exactly tangential and has minimized local distortion [23]. After that, we cluster the projected triangles to atlases via connected components,

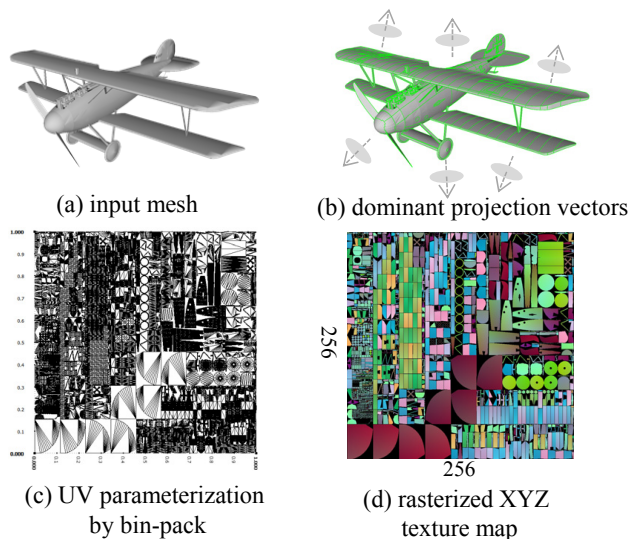


Figure 1: Given an input mesh (a), we compute its dominant projections (b), and then pack all atlases in one UV map (c). The specific feature (such as vertex positions) is rasterized in the texture map (d).

---

#### Algorithm 1 virtual texture map generation for pure mesh

---

```

// input: the mesh  $\mathcal{M} = \{V, T\}$ 
// params:  $\tau_{\text{angle}} = 30^\circ$ 
// step 1: finding dominant projection vectors  $\mathbf{P}$ 
 $\mathbf{P} = \{\}$  // the projection vector set
sort  $T$  based on triangle areas in descending order
for triangle  $t_i$  in  $T$  and  $\text{visited}(t_i)$  is false:
     $\mathbf{C} = \{t_i\}$  // the candidate set
    for triangle  $t_j$  in  $T$ ,  $\text{visited}(t_j)$  is false and  $i \neq j$ :
        if  $\text{normal}(t_i) \cdot \text{normal}(t_j) > \cos(\tau_{\text{angle}})$ :
             $\mathbf{C} = \mathbf{C} \cup \{t_j\}$ 
     $P = (\sum_{t \in \mathbf{C}} \text{area}(t) \cdot \text{normal}(t)) / \sum_{t \in \mathbf{C}} \text{area}(t)$ 
     $\mathbf{P} = \mathbf{P} \cup \{P\}$ 
     $\text{visited}(\forall t \in \mathbf{C}) = \text{true}$ 
// step 2: assign triangles to the best projection vector
for triangle  $t_i$  in  $T$ :
    for projection vector  $P_j$  in  $\mathbf{P}$ :
        if  $(P_j \cdot \text{normal}(t_i) > \text{view}(t_i) \cdot \text{normal}(t_i))$ :
             $\text{view}(t_i) = P_j$ 
// step 3: create UV projection and bin-packing
orthogonally project  $T$  from 3D to 2D using  $\text{view}(t_i)$ 
cluster connected components to atlases, and pack them
into a squared UV map them using [15]
// step 4: rasterize the texture map for network input
rasterize the texture map  $\mathbf{T}$  with resolution  $H \times W$  and
the choice of a  $C$ -channel feature vector
// output: texture map  $\mathbf{T}$ 

```

---



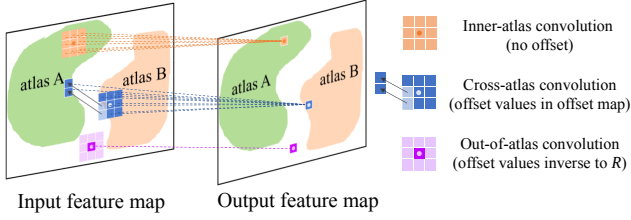


Figure 2: Illustration of the cross-atlas convolution. There are three situations when applying convolution over three regions.

and pack them into one squared map using the bin-packing algorithm [15]. With this UV parameterization, we can rasterize the geometric feature of the mesh to a  $H \times W$  resolution texture map  $\mathbf{T}$ . Note that the pixel in texture map  $\mathbf{T}(x, y) = [f_1, f_2, \dots, f_c]^T$  is a  $C$ -dimensional feature vector, instead of RGB values of the standard definition of “texture maps”. The choice of geometric feature could be intrinsic features (e.g., curvatures, heat kernel signatures) or extrinsic properties (e.g., spatial coordinates, normals), or even concatenate multiple of them, depending on the specific task. The output texture map is a  $H \times W \times C$  tensor, where  $H \times W$  is the spatial resolution and  $C$  is the feature vector channel.

If the mesh comes with textures, we still generate our parameterization using Algorithm 1. The RGB color in original texture maps is an additional feature that can be concatenated to the geometric feature vector. We do not use the parameterization given in the original textured mesh as we need to ensure two criteria of parameterization. 1) It should be *area-preserving*, which means the areas of triangles in meshes and their projected areas in 2D maps are proportional. It ensures the receptive field of 2D convolution over the texture map should correspond to equal geodesic area over the mesh surface. This importance is also mentioned in [30]. 2) It should be *rotation-aligned*: the rotation of atlas in the original texture map could be arbitrary, making the later convolution suffer from rotation ambiguities. In our generated texture map, the rotation of each atlas is aligned with the negative Z-axis (usually the gravity direction) of the mesh coordinates, such that the visual content of atlas is upright.

Unlike natural images, the atlas in texture map is discontinuous and locates randomly. It is not visually meaningful, and standard CNN approaches would not take effect on such input. To bridge the atlases and apply convolution across them, we also need to generate the offset map of size  $(H \times W \times 2k^2)$ ,  $k$  is the kernel size, which encodes the neighborhood information between atlas boundaries. We will describe how to create it in the next section together with the cross-atlas convolution.

### 3.2. Cross-atlas convolution via kernel offsets

The standard 2D convolution computes the output feature map  $\mathbf{F}_o$  via the weighted sum of a  $k$  size regular patch over every pixel  $\mathbf{p} = (x, y)$  at the input feature map  $\mathbf{F}_i$ :

$$\mathbf{F}_o(\mathbf{p}) = \sum_{\mathbf{p}_n \in \mathcal{R}} \mathbf{F}_i(\mathbf{p} + \mathbf{p}_n) \cdot g(\mathbf{p}_n), \quad (1)$$

where  $g(\cdot)$  is the kernel weight and  $\mathcal{R} = \{\mathbf{p}_n\} = \{(x, y) : -\frac{k-1}{2} \leq x, y \leq \frac{k-1}{2}\}$  enumerates neighboring locations of the center pixel and indicates the receptive field.

This equation holds when the local neighborhood is simply the  $k \times k$  block in natural images. For texture maps where atlases are isolated, two neighboring surface points on mesh can lie on two separated atlases. To recover the original mesh geodesic neighborhood, the 2D convolution should be able to apply across atlases.

To this end, when rasterizing the texture maps in Section 3.1, we encode the atlas connectivity information by generating the corresponding offset map  $\mathcal{R}_{offset}$  with the same spatial resolution  $H \times W$  and channel length  $2k^2$ , where each pixel  $\mathbf{p}$  in the offset map has a  $k \times k$  block indicating the real-valued offsets of x-axis and y-axis:

$$\mathcal{R}_{offset}(\mathbf{p}) = \{(\Delta x, \Delta y)\} \quad (2)$$

$$= \{\Delta \mathbf{p}_n\}, \quad |\mathcal{R}_{offset}(\mathbf{p})| = k^2. \quad (3)$$

These offset values are augmented to the standard neighboring locations  $\mathcal{R} = \{\mathbf{p}_n\}$  and redirect to the pixel corresponding to mesh geodesic neighboring point, i.e.,  $\mathcal{R} + \mathcal{R}_{offset}$  indicates the neighborhood, and  $\mathbf{p} + \mathbf{p}_n + \Delta \mathbf{p}_n$  is the geodesic neighboring location of  $\mathbf{p}$ . The original equation in Eq. 1 becomes

$$\mathbf{F}_o(\mathbf{p}) = \sum_{\substack{\mathbf{p}_n \in \mathcal{R} \\ \Delta \mathbf{p}_n \in \mathcal{R}_{offset}}} \mathbf{F}_i(\mathbf{p} + \mathbf{p}_n + \Delta \mathbf{p}_n) \cdot g(\mathbf{p}_n). \quad (4)$$

As illustrated in Fig. 2, we classify the regions in the texture map into three types, denoted by pixels  $\mathbf{p}_{in}, \mathbf{p}_{out}, \mathbf{p}_{across}$ . 1) When the convolution is applied over the inner-atlas region, the standard pixel neighborhood is corresponding to the mesh geodesic neighborhood, and thus no offset should be added:  $\mathcal{R}_{offset}(\mathbf{p}_{in}) = \{(0, 0)\}$ . 2) When applying over the out-of-atlas region, this pixel value is invalid and should be kept isolated in order not to contaminate other pixels. We use offset values inverse to  $\mathcal{R}$ , i.e.,  $\mathcal{R}_{offset}(\mathbf{p}_{out}) = -\mathcal{R}$  and  $\mathbf{p}_{out} + \mathbf{p}_n + \Delta \mathbf{p}_n = \mathbf{p}_{out}$ . This ensures  $\mathbf{F}_o(\mathbf{p}_{out}) = \mathbf{F}_i(\mathbf{p}_{out}) = 0$  throughout the network. 3) When the standard  $\mathcal{R}$  is just across the border of an atlas, we add the precomputed offset values to its standard locations, so  $\mathbf{p}_{across} + \mathbf{p}_n + \Delta \mathbf{p}_n$  should just locate at the true mesh geodesic neighborhood. Therefore, only pixels  $\mathbf{p}_{across}$  need to compute their offset values. The pixels

$\mathbf{p}_{across}$  are determined by the kernel size: if we place a  $k$ -size filter kernel within the atlas at  $\mathbf{p}$  and there are some pixels of this kernel locate out of the atlas, then  $\mathbf{p} = \mathbf{p}_{across}$ . For a  $3 \times 3$  kernel as an example, we compute the offsets for 1-ring atlas boundary pixels.

To compute the offset value for a center pixel  $\mathbf{p}$  regarding its neighbor pixel  $\mathbf{p} + \mathbf{p}_n$  lying out of atlas, we rasterize the vertex coordinates to a map with the original resolution  $H \times W$ , and thus can instantly query via this map the pixel  $\mathbf{p}$  corresponds to the 3D point  $\mathbf{X}$  on mesh surface. Then we use the Fast Marching [33] to search the geodesic neighbor point of  $\mathbf{X}$  along  $\vec{\mathbf{p}}_n$  direction, yielding the point  $\mathbf{X}'$ . The  $\mathbf{X}'$  finds its corresponding texture coordinates  $\mathbf{p}'$  on the texture map. Finally  $\Delta\mathbf{p} = \mathbf{p}' - \mathbf{p} - \mathbf{p}_n$  is the offset value.

Note that unlike standard convolution on natural images can add paddings to image boundaries so as to keep the same dimension between input and output feature maps, the cross-atlas convolution has no “image boundaries” – if the mesh surface is water tight, every pixel in texture map can find its neighborhood. If it is an open mesh and the pixel exactly corresponds to the mesh boundary that finds no neighborhood, it simply fills zero value (analogous to the zero padding effect). In all, the spatial resolution of feature map can only be decreased by either strides or pooling.

**Deconvolution** In semantic segmentation tasks where the feature map is finally up-sampled to the input resolution, deconvolution is a popular approach [18]. Essentially, deconvolution (also dubbed transposed convolution) can be decomposed into two operations: 1) scattering the pixels from the sparse feature map to a dense feature map with strides and 2) applying convolution over this map. We can simply replace the convolution in 2) with the cross-atlas version if applied on texture maps.

**Pooling** The standard pooling operation downsamples the input feature map and reduce the number of parameters. The pooling type can be *max*, *average* or *sum* operation.

$$\mathbf{F}_o(\mathbf{p}_o) = \text{pool}_{k \times k}(\mathbf{F}_i(\mathbf{p}), \mathbf{F}_i(\mathbf{p} + \mathbf{p}_n), \dots), \forall \mathbf{p}_n \in \mathcal{R} \quad (5)$$

where  $\mathbf{p}_o$  is the location at the lower spatial resolution output feature map. For  $k = 2$  as an example,  $\mathbf{p}_o \cdot 2 = \mathbf{p}$ . It computes the pooling value from every  $2 \times 2$  block.

In cross-atlas pooling, the behavior is similar as cross-atlas convolution by replacing the standard image pixel neighborhood with the mesh geodesic neighborhood:

$$\mathbf{F}_o(\mathbf{p}_o) = \text{pool}_{k \times k}(\mathbf{F}_i(\mathbf{p}), \mathbf{F}_i(\mathbf{p} + \mathbf{p}_n + \Delta\mathbf{p}_n), \dots), \quad \forall \mathbf{p}_n \in \mathcal{R}, \Delta\mathbf{p}_n \in \mathcal{R}_{offset}. \quad (6)$$

**Hierarchies** For a specific  $H \times W$  dimension feature map and the kernel size  $k$ , the corresponding offset map is unique. In the CNN pipeline, the spatial dimension of the

feature map keeps changing throughout the network. Therefore, the corresponding offset maps for all possible feature map dimensions need to be computed beforehand. Generally, we compute the pyramid of offset maps from the original size, where each lower level is half in width and height of the upper level. Note that an offset map with larger kernel size can be reused in the operation with smaller kernel. *e.g.*, we can take the central  $3 \times 3$  from  $5 \times 5$  offset locations.

There are some similarities between our design and the deformable convolution [7] as we both leverage offset values to the convolution. However, some crucial points are different: 1) our offset is pre-computed by the atlas neighborhood, while offset values in [7] are all trainable. Their method introduces more parameters to be learned, while ours has no extra parameters. This allows us to train on natural images with standard CNN, and test on textured meshes with cross-atlas convolution (see also Section 5.3). 2) The objective is different: their work aims to find a more adaptive receptive field. Our work recovers the mesh geodesic neighborhood. 3) They propose the region-of-interest (RoI) pooling for object detection, whereas our work defines the cross-atlas version of deconvolution and pooling.

### 3.3. Network architecture

We have integrated the cross-atlas convolution into classification and semantic segmentation architectures.

**Classification** A generic classification network inputs an image or a feature map, which goes through multiple layers of convolution and pooling. Over these layers, the spatial resolution is decreased by pooling or strides, and the number of feature channel is increased by convolution. Finally it is flattened to a 1D global feature vector, followed by fully connected (FC) layers and softmax, yielding the class label.

To apply classification on (textured) mesh, we first convert the mesh to the  $H_1 \times W_1 \times C_1$  input feature map (Section 3.1) and corresponding offset maps. Then, we replace the standard convolution and pooling with our cross-atlas versions. After bypassing  $n$  layers of cross-atlas convolutions, it obtains a feature map with dimensions  $H_n \times W_n \times C_n$ . At this moment, we cannot simply flatten it to a 1D feature vector like the standard network does, because the spatial location of each pixel in this feature map is completely permutable if we pack the atlases in a different way in the input texture map. Inspired by PointNet [24], we regard each pixel in the  $H_n \times W_n \times C_n$  feature map is a permutable “point”. We reshape it to  $(H_n \times W_n) \times C_n \times 1$ , *i.e.*, each  $C_n$  pixel vector is expanded to one row. Then we apply multi-layer perceptrons (MLP) and row-wise max pooling to obtain a  $m$ -channel 1D feature vector. Finally fully connected layers is applied and outputs the  $n_{class}$ -channel 1D feature vector indicating the probability of each class. The specific number of layers and the dimension of feature map vary depending on the complexity of each task.

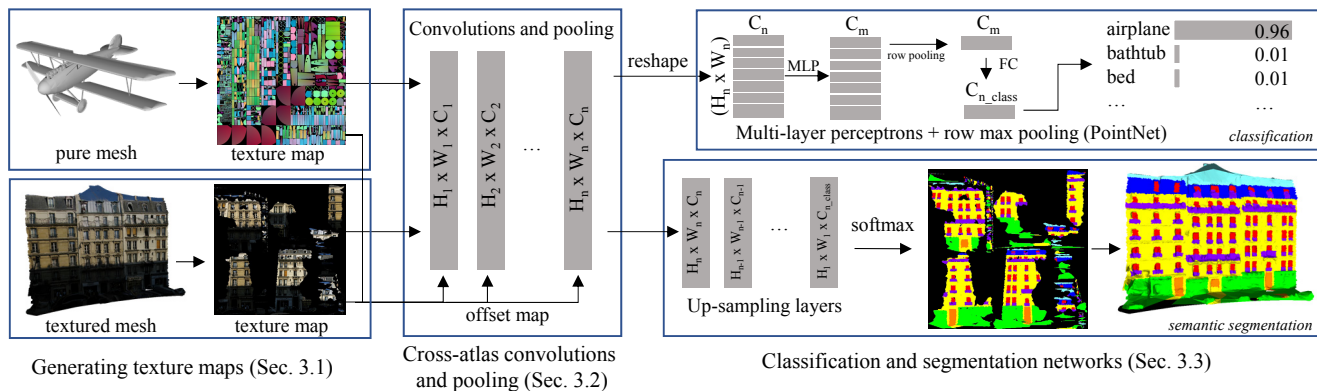


Figure 3: The classification and segmentation network architecture of our method.

**Segmentation** A segmentation task is a pixel-wise classification task. Its former part is similar to classification which yields the  $H_n \times W_n \times C_n$  feature map via several layers of convolution and pooling. Unlike classification which extracts the global feature in the later modules, it up-samples the feature map to an original resolution annotation map  $H_1 \times W_1 \times C_{n_{class}}$ . Here, we follow the deconvolution operations [18] (or dubbed transposed convolution) as the up-sampling layers, and replace their convolution with our cross-atlas version.

In the experiment section, we evaluate these network architects on classification and segmentation tasks. The concrete number of layers used in each task will be described.

#### 4. Understanding the cross-atlas convolution

In this section, we discuss three important properties of the proposed method. 1) It is robust to irregular meshes. 2) It is invariant to parameterization. 3) The receptive field follows the mesh geodesic distance.

**Robust to irregular meshes** By the design of mesh, a similar shape can be composed by very different meshing structures, but it is not expected different meshing would cause inconsistent results in mesh learning. To tackle this, the texture map rasterization in our first step can be deemed as the procedure of sampling the surface points to pixels. If we connect every 4-neighborhood of pixels to edges, it reconstructs a regular mesh as shown in Fig. 4. Therefore, applying convolution over the texture map is equivalent to applying convolution over this mesh. When a lower resolution ( $256 \times 256$ ) texture map is used, thin structures such as bed legs are missing (Fig. 4(b)), whereas higher resolution ( $1024 \times 1024$ ) involves more geometric details (Fig. 4(d)).

**Invariant to parameterization** When generating the texture map (Section 3.1), the locations of texture atlas are unpredictable and even random every time we generate it.

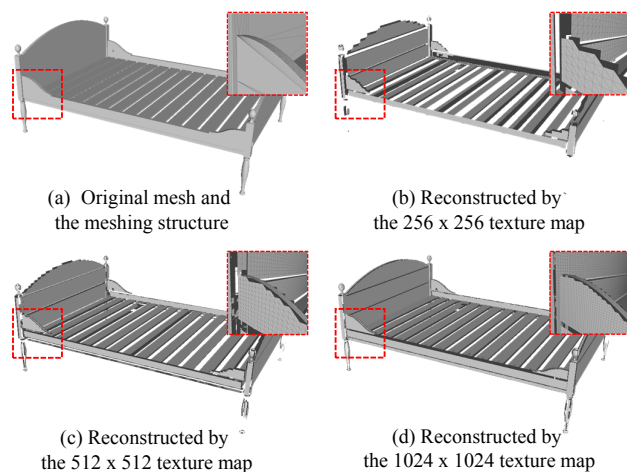


Figure 4: The original mesh (a) and reconstructed meshes from texture maps with different resolution (b)(c)(d). Our method is invariant to any irregular mesh input: applying convolution over the texture map is equivalent to applying over the vertices of these reconstructed meshes.

The former part of the network pipeline includes multiple layers of cross-atlas convolution and pooling. These operations are translational equivalent, *i.e.*, if we have the correct offset map, the receptive field of convolution always follows the mesh geodesic range, no matter how the atlas is parameterized in 2D texture map. Then, in the later part of classification, we use a combination of MLP, max pooling and FC (same as PointNet [24]) to extract global feature, which is invariant to the spatial locations of pixel-wise feature vector in the last feature map ( $H_n \times W_n \times C_n$ ) of “cross-atlas convolution” module. Regarding segmentation task, the later part of up-sampling layers is deconvolution operation, which enjoys the same translational equivalent property as convolution.

Note that the parameterization here only refers to arbitrary translations, which indicate how atlases are packed in



the texture map. For rotations, we have already aligned all atlases with Z-axis to disambiguate X-axis and Y-axis rotation varieties. Although there still exists ambiguities with regards to the in-plane rotation of triangles perpendicular to Z-axis, we can augment the training data by rotating along Z-axis as most previous works do [21, 30, 24].

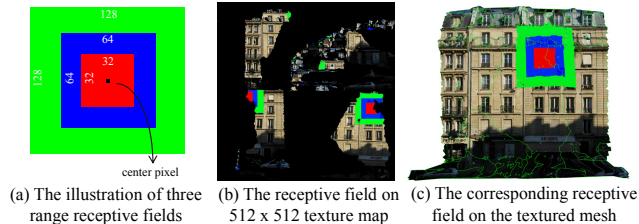


Figure 5: Illustration of the receptive field. If we dilate the receptive field from a center pixel on the texture map and redirect when reaching to atlas boundaries, the field would be separated in several atlases (b). Its corresponding field on textured mesh is approximately the geodesic field.

**Receptive fields** Although the convolution is applied on 2D texture maps, its receptive field follows the mesh geodesic distance. This is done by using the offset map to redirect the locations. Fig. 5 illustrates this behavior: we color-code the receptive fields of  $32 \times 32$ ,  $64 \times 64$ ,  $128 \times 128$  of a center pixel in red, blue and green respectively. We can see in the texture map, as we dilate the field and redirect to the offset location when reaching to atlas boundaries, the field is actually separated in several atlases. On the contrary, its corresponding textured mesh shows the correct block-wise receptive field over its surface. Noted that the receptive field on the mesh is not strictly following the geodesic distance due to the distortion in projections, but it is a good approximation given the distortion is controllably small.

## 5. Experiments

We implement the texture map generation (Section 3.1) in C++, and the network (Section 3.3) using Tensorflow. Our method is evaluated on three benchmarks. First we conduct ablation studies on the MNIST dataset, and then we test the classification performance on ModelNet [38] and segmentation performance on Rue Monge 2014 dataset [28].

### 5.1. MNIST Textured meshes

The original MNIST dataset contains 60,000 training images and 10,000 testing images of 10 handwritten digits at  $28 \times 28$  resolution (Fig. 6(a)). We create the textured mesh version by mapping the digit image to a planar mesh, which is triangulated from 1,000 randomly distributed points (Fig. 6(b)). Then, we segment the mesh and pack atlases in the texture map (Fig. 6(c)). The texture map and its corresponding offset map are fed into our network to train a classifier.

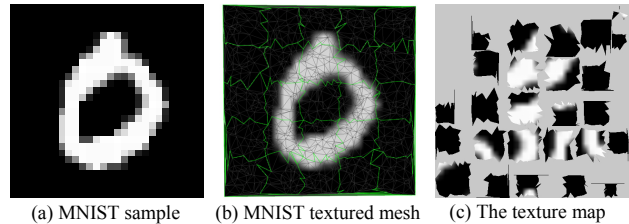


Figure 6: The MNIST data sample (a) is texture mapped to a mesh (b), and its texture map is not visually recognizable as in (c).

**Ablation study** To better validate the effectiveness of each component, we start with the standard LeNet5, and then replace with our components one by one. Table 1 shows the comparisons of using different modules and dataset.

conv. layers	FC layers	dataset	acc.
standard conv.	standard FC	digit image	99.2%
		texture maps	32.5%
standard conv.	MLP+max pooling+FC	digit image	98.3%
		texture maps	91.6%
cross-atlas conv.	MLP+max pooling+FC	digit image	98.3%
		texture maps	97.2%

Table 1: The testing accuracy of combinations of different modules and datasets. Our method achieves 97.2% accuracy on the textured mesh version of MNIST, which supports the effectiveness of proposed method.

The standard LeNet5 consists of four convolutional layers and three fully connected (FC), achieving 99.2% on the original MNIST dataset, and 32.5% on texture maps, which indicates that the network cannot learn a good model to distinguish texture maps. If we replace the standard FC layers with the MLP+max pooling+FC modules, the accuracy on digit image drops by 0.9% on original dataset and increases to 91.6% on texture maps. The dropped performance in the standard images is caused by the loss of spatial information our FC layers, but it dramatically improves the performance by  $\sim 60\%$  on texture maps as it achieves translational invariance of atlases in the texture map. If we further integrates the cross-atlas convolution and pooling in the network, the accuracy on texture maps increases to 97.2%.

### 5.2. ModelNet classification

We evaluate our approach for 3D shape classification task on the two versions of the large scale Princeton ModelNet dataset [38]: ModelNet40 and ModelNet10, which consist of 40 and 10 classes respectively. We follow the same experiment setting as in [38]. The vertex coordinates are rasterized to the texture map at  $256 \times 256$  resolution for the network input. As our texture bin-packing algorithm is randomized, we generate the input texture map by running multiple times to augment the training data.

Method	input	ModelNet40 accuracy	ModelNet10 accuracy
RotationNet [11]	image	97.37%	98.46%
VoxNet [21]	volume	83%	92%
PointNet++ [26]	point	91.9%	-
SHR [12]	mesh	68.2%	79.9%
GeometryImage[30]	mesh	83.9%	88.4%
Ours	mesh	87.5%	91.2%

Table 2: The overall classification accuracies of the multi-view image, volume, point and mesh representations.

Table 2 shows the classification accuracy in testing. We have listed representative state-of-the-art methods of using different geometry representations, namely multi-view images, volumes, points and meshes. Our method achieve better results than the other two mesh-based methods [12, 30], while it is overall not as competitive as multi-view image projection or point-based methods. We perceive the CAD mesh has a signification problem that hinders the performance of mesh-based methods: some structure in the mesh model should have been topologically connected, but in fact they are just overlaid together. This makes the geodesic receptive field erroneous. On the contrary, multi-view image projection or point-based method can avert the problem. This also infers our method works better on the “natural” meshes created by laser-scan or photogrametry, rather than the “artificial” meshes where the topology of the mesh is not reliable.

### 5.3. Ruemonge2014 segmentation

Method	triangle accuracy	class avg. IoU
Riemenschneider [28]	-	41.92%
Gadde [9]	-	63.7%
Ours (texture, w/o fusion)	72.28%	65.24%
Ours (texture, w/ fusion)	84.90%	75.34%
Ours (image, w/ fusion)	76.02%	72.38%
Ours (texture+image, w/ fusion)	85.63%	75.67%

Table 3: The evaluation of mesh-based semantic segmentation on the Ruemonge2014 dataset [28]

We evaluate the semantic segmentation performance on the Ruemonge2014 dataset [28], which consists of 428 high resolution images capturing the facade along the street, as well as registered camera poses and reconstructed meshes by multi-view stereo. The images and mesh come with ground truth semantic labels of seven classes, namely the door, shop, balcony, window, wall, sky and rooftop. The dataset is separated into training samples and testing samples. To evaluate, the class-averaged intersection over union (IoU) are per-triangle label accuracy are used.

To generate the textured mesh, we run the multi-view

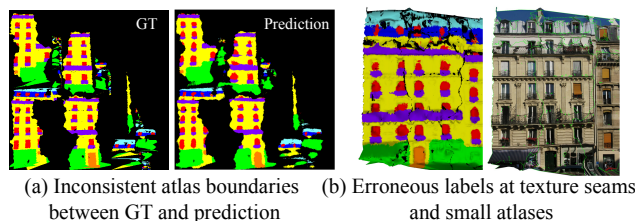


Figure 8: The problem of semantic segmentation on texture maps: the atlas boundaries are not always consistent with ground truth (a), leading to erroneous labels at texture seams (b).

texturing algorithm [34] using the given undistorted images and camera poses. Then, we segment the training mesh part into 100 overlapped mesh segments, each consists of a  $512 \times 512$  texture map. We integrate our cross-atlas convolution modules into the fully convolutional networks (FCN) [18] with VGG19 [29]. This network exactly corresponds to our generic segmentation architecture in Fig. 3: the VGG part corresponds to the cross-atlas convolution. The up-sampling layers are identical to the ones used in the FCN. We replace all  $3 \times 3$  convolutions and  $2 \times 2$  pooling in VGG and the deconvolution in FCN with ours cross-atlas version and apply on the texture maps with offset maps. Here, the RGB+height values are used in the texture map channel. The training parameters details are in the supplementary material.

Fig. 8 shows one testing result. We notice that the predicted annotation map does not have exactly the same atlas boundaries as the ground truth, and some small atlases are even filtered out (Fig. 8(a)). This is potentially due to the lost of spatial location information in up-sampling layers in FCN [18]: in natural images, the annotation boundaries may have “over-rounded” artifacts, whereas in texture maps, the atlas may slightly dilate or erode. This issue leads to erroneous labels at texture seams when mapping the annotation map to the mesh (Fig. 8(b)).

**Fusion in testing stage** Inspired by the multi-scale trick used in the testing stage of semantic segmentation methods, we have tried conducting the testing on individual overlapped parts, and fuse them afterwards: each triangle label is finally determined by the label of majority pixels. This improves the result significantly, from 65.24% to 77.34% IoU as shown in Table 3. Overall, our method surpass the second place in the Ruemonge2014 challenge benchmark [28] by a large margin. Fig. 7 shows the qualitative results.

**Train on images, test on textured meshes** With cross-atlas convolution, the semantic segmentation network can be trained and tested on texture maps with corresponding offset maps. One may wonder if we can also train on normal



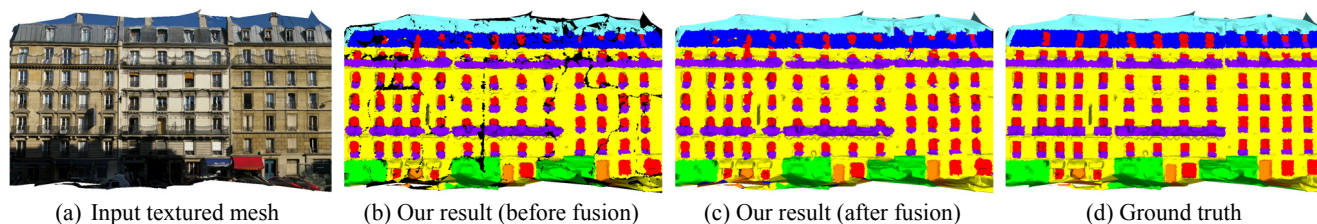


Figure 7: The qualitative comparison of semantic segmentation results on the RueMonge2014 dataset [28].

images and test on texture maps – we can regard the image is a one-atlas “texture map” with zero offsets. To validate, we take street view images and their corresponding ground truth given in the dataset for training. The testing accuracy turns out decent, achieving 76.02% triangle label accuracy and 72.38% IoU. If we combine the normal images and texture maps for training, the result is by 0.33% marginally better than only training on texture maps. This shows our method is flexible in terms of the training data – it can be trained on normal images and test on texture maps of meshes with the corresponding content in original images.

## 6. Conclusion

We have proposed a parameterization invariant approach to textured mesh learning. The key to this method is the cross-atlas convolution which recovers the mesh geodesic receptive field although it actually convolves on 2D domain. Our work unlocks the possibility for direct classification and semantic segmentation on textured meshes via their texture maps, whereas in previous methods this only can be done by multi-view image projection or point cloud sampling from meshes.

**Limitations and future works** The biggest limitation of the proposed method is the texture map requires the rotation of atlases should be aligned to upright, so the convolution is not suffer from rotation variances of two axes. For the plane perpendicular to Z-axis, there’s still in-plane rotation ambiguity, which we resolve it via augmenting the training data. Some related works use polar convolution [20, 2, 22] or pooling after multi-directional convolution [23] to achieve fully rotational invariance, as their shape correspondence task is more sensitive to geometries. We do not use them in our current framework but it could be a potential improvement [37]. Besides, object detection over textured meshes is another potential task, which is well-addressed in 2D images but rarely targets to 3D models.

## References

- [1] F. Bogo, J. Romero, M. Loper, and M. J. Black. FAUST: Dataset and evaluation for 3D mesh registration. In *Proceed-*

- ings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, Piscataway, NJ, USA, June 2014. IEEE. 1
- [2] D. Boscaini, J. Masci, E. Rodolà, and M. Bronstein. Learning shape correspondence with anisotropic convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 3189–3197, 2016. 1, 3, 9
- [3] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017. 1, 2
- [4] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013. 1, 2
- [5] A. Dai, A. X. Chang, M. Savva, M. Halber, T. A. Funkhouser, and M. Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, volume 2, page 10, 2017. 1, 2
- [6] A. Dai and M. Nießner. 3dmv: Joint 3d-multi-view prediction for 3d semantic scene segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 1, 2
- [7] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei. Deformable convolutional networks. *CoRR, abs/1703.06211*, 1(2):3, 2017. 5
- [8] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems*, pages 3844–3852, 2016. 1, 2
- [9] R. Gadde, V. Jampani, R. Marlet, and P. V. Gehler. Efficient 2d and 3d facade segmentation using auto-context. *IEEE transactions on pattern analysis and machine intelligence*, 40(5):1273–1280, 2018. 8
- [10] E. Kalogerakis, M. Averkiou, S. Maji, and S. Chaudhuri. 3d shape segmentation with projective convolutional networks. In *Proc. CVPR*, volume 1, page 8, 2017. 1, 2
- [11] A. Kanezaki, Y. Matsushita, and Y. Nishida. Rotationnet: Joint object categorization and pose estimation using multi-views from unsupervised viewpoints. In *Proc. 2018 IEEE Conf. on Computer Vision and Pattern Recognition*, pages 5010–5019, 2018. 1, 2, 8
- [12] M. Kazhdan, T. Funkhouser, and S. Rusinkiewicz. Rotation invariant spherical harmonic representation of 3 d shape descriptors. In *Symposium on geometry processing*, volume 6, pages 156–164, 2003. 8

- [13] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016. 1
- [14] R. Klokov and V. Lempitsky. Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 863–872. IEEE, 2017. 2
- [15] R. E. Korf. A new algorithm for optimal bin packing. In *AAAI/IAAI*, pages 731–736, 2002. 3, 4
- [16] I. Kostrikov, Z. Jiang, D. Panozzo, D. Zorin, and B. Joan. Surface networks. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018*, 2018. 1, 2
- [17] Y. Li, R. Bu, M. Sun, and B. Chen. Pointcnn. *arXiv preprint arXiv:1801.07791*, 2018. 2
- [18] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 5, 6, 8
- [19] H. Maron, M. Galun, N. Aigerman, M. Trope, N. Dym, E. Yumer, V. G. Kim, and Y. Lipman. Convolutional neural networks on surfaces via seamless toric covers. *ACM Trans. Graph*, 36(4):71, 2017. 2, 3
- [20] J. Masci, D. Boscaini, M. Bronstein, and P. Vandergheynst. Geodesic convolutional neural networks on riemannian manifolds. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 37–45, 2015. 1, 3, 9
- [21] D. Maturana and S. Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*, pages 922–928. IEEE, 2015. 1, 2, 7, 8
- [22] F. Monti, D. Boscaini, J. Masci, E. Rodola, J. Svoboda, and M. M. Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Proc. CVPR*, volume 1, page 3, 2017. 1, 3, 9
- [23] H. Pan, S. Liu, Y. Liu, and X. Tong. Convolutional neural networks on 3d surfaces using parallel frames. *arXiv preprint arXiv:1808.04952*, 2018. 3, 9
- [24] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 1(2):4, 2017. 1, 2, 5, 6, 7
- [25] C. R. Qi, H. Su, M. Nießner, A. Dai, M. Yan, and L. J. Guibas. Volumetric and multi-view cnns for object classification on 3d data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5648–5656, 2016. 1, 2
- [26] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems*, pages 5099–5108, 2017. 1, 2, 8
- [27] G. Riegler, A. O. Ulusoy, and A. Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 3, 2017. 2
- [28] H. Riemenschneider, A. Bódis-Szomorú, J. Weissenberg, and L. Van Gool. Learning where to classify in multi-view semantic segmentation. In *European Conference on Computer Vision*, pages 516–532. Springer, 2014. 7, 8
- [29] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 8
- [30] A. Sinha, J. Bai, and K. Ramani. Deep learning 3d shape surfaces using geometry images. In *European Conference on Computer Vision*, pages 223–240. Springer, 2016. 2, 3, 4, 7, 8
- [31] H. Su, V. Jampani, D. Sun, S. Maji, E. Kalogerakis, M.-H. Yang, and J. Kautz. SPLATNet: Sparse lattice networks for point cloud processing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2530–2539, 2018. 1, 2
- [32] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015. 1, 2
- [33] V. Surazhsky, T. Surazhsky, D. Kirsanov, S. J. Gortler, and H. Hoppe. Fast exact and approximate geodesics on meshes. In *ACM transactions on graphics (TOG)*, volume 24, pages 553–560. Acn, 2005. 5
- [34] M. Waechter, N. Moehle, and M. Goesele. Let there be color! large-scale texturing of 3d reconstructions. In *European Conference on Computer Vision*, pages 836–850. Springer, 2014. 8
- [35] P.-S. Wang, Y. Liu, Y.-X. Guo, C.-Y. Sun, and X. Tong. O-cnn: Octree-based convolutional neural networks for 3d shape analysis. *ACM Transactions on Graphics (TOG)*, 36(4):72, 2017. 2
- [36] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon. Dynamic graph cnn for learning on point clouds. *arXiv preprint arXiv:1801.07829*, 2018. 2
- [37] D. E. Worrall, S. J. Garbin, D. Turmukhambetov, and G. J. Brostow. Harmonic networks: Deep translation and rotation equivariance. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, volume 2, 2017. 9
- [38] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015. 1, 2, 3, 7
- [39] L. Yi, H. Su, X. Guo, and L. J. Guibas. Syncspecnn: Synchronized spectral cnn for 3d shape segmentation. In *CVPR*, pages 6584–6592, 2017. 1, 2