

# Pix2Surf: Learning Parametric 3D Surface Models of Objects from Images

Jiahui Lei<sup>1</sup>   Srinath Sridhar<sup>2</sup>   Paul Guerrero<sup>3</sup>   Minhyuk Sung<sup>3</sup>  
 Niloy Mitra<sup>4,3</sup>   Leonidas J. Guibas<sup>2</sup>

<sup>1</sup>Zhejiang University   <sup>2</sup>Stanford University   <sup>3</sup>Adobe Research  
<sup>4</sup>University College London

🌐 <https://geometry.stanford.edu/projects/pix2surf>

**Abstract.** We investigate the problem of learning to generate 3D parametric surface representations for novel object instances, as seen from one or more views. Previous work on learning shape reconstruction from multiple views uses discrete representations such as point clouds or voxels, while continuous surface generation approaches lack multi-view consistency. We address these issues by designing neural networks capable of generating *high-quality parametric 3D surfaces* which are also *consistent* between views. Furthermore, the generated 3D surfaces preserve accurate image pixel to 3D surface point *correspondences*, allowing us to lift texture information to reconstruct shapes with rich geometry *and* appearance. Our method is supervised and trained on a public dataset of shapes from common object categories. Quantitative results indicate that our method significantly outperforms previous work, while qualitative results demonstrate the high quality of our reconstructions.

**Keywords:** 3D reconstruction, multi-view, single-view, parametrization

## 1 Introduction

Reconstructing the 3D shape of an object from one or more views is an important problem with applications in 3D scene understanding, robotic navigation or manipulation, and content creation. Even with multi-view images, the problem can be challenging when camera baselines are large, or when lighting and occlusions are inconsistent across the views. Recent developments in supervised deep learning have demonstrated the potential to overcome these challenges.

Ideally, a multi-view surface reconstruction algorithm should have the following desirable 3C properties: surface continuity, multi-view consistency and 2D-3D correspondence. First, it should be able to reconstruct high-quality shapes that can be readily used in downstream applications. While much progress has been made in learning shape representations such as point clouds [11,25,39,17], volumetric grids [9,41,42], and meshes [44,47], their geometric quality is limited by the discrete nature of the underlying representation. Therefore, representations such as implicit functions [34,37,7], and *UV* surface parametrizations [15,10] are



**Fig. 1.** Pix2Surf learns to generate a continuous parametric 3D surface of an object seen in one or more views. Given a single image, we can reconstruct a continuous partial 3D shape (top row). When multiple views are available, we can aggregate the views to form a set of consistent 3D surfaces (bottom row). Our reconstructions preserve 2D pixel to 3D shape correspondence that allows the transport of textures, even from real images (last column).

preferable, since they can represent a **continuous surface** at arbitrary resolution. Second, the algorithm should be able to reconstruct objects from a sparse set of views while ensuring that the combined shape is **consistent** across the views. Recent approaches exploit geometric constraints to solve this problem but require additional supervision through knowledge of the exact camera geometry [5]. Finally, the algorithm should provide accurate **correspondences** between 2D pixels and points on the 3D shape, so as to accurately transport object properties (e.g., texture) directly from 2D and support aggregation across views. While some extant methods satisfy a subset of these properties, we currently lack any method that has all of them.

In this paper, we present **Pix2Surf**, a method that learns to reconstruct *continuous* and *consistent* 3D surface from single or multiple views of novel object instances, while preserving accurate 2D–3D *correspondences*. We build upon recent work on category-specific shape reconstruction using *Normalized Object Coordinate Space (NOCS)* [43,39], which reconstructs the 3D point cloud as a *NOCS map* – an object-centered depth map – in a canonical space that is in accurate correspondence with image pixels. Importantly, NOCS maps do not require knowledge of camera geometry. However, these maps do not directly encode the underlying surface of the object. In this paper, we present a method that incorporates a representation of the underlying surface by predicting a continuous parametrization that maps a learned *UV parameter space* to 3D NOCS coordinates, similar in spirit to AtlasNet [15]. Unlike AtlasNet, however, our approach also provides accurate 2D–3D correspondences and an emergent learned chart that can be used to texture the object directly from the input image.

When multiple views of an object are available, we also present a version of Pix2Surf that is capable of reconstructing an object by predicting an atlas, i.e., view-specific charts assembled to form the final shape. While in the NOCS approach [39] individual views can also be directly aggregated since they live in

the same canonical space, this naïve approach can lead to discontinuities at view boundaries. Instead, for view-consistent reconstruction, we aggregate multiple views at the feature level and explicitly enforce consistency during training.

Extensive experiments and comparisons with previous work show that Pix2Surf is capable of reconstructing high-quality shapes that are consistent within and across views. In terms of reconstruction error, we outperform state-of-the-art methods while maintaining the 3C properties. Furthermore, accurate 2D–3D correspondences allow us to texture the reconstructed shape with rich color information as shown in Fig. 1. In summary, the primary contributions of our work are:

- a method to generate a set of **continuous** parametric 3D surfaces representing the shape of a novel object observed from single or multiple views;
- the unsupervised extraction of a learned  $UV$  parametrization that retains accurate 2D to 3D surface point **correspondences**, allowing lifting of texture information from the input image; and
- a method to **consistently** aggregate such parametrizations across different views, using multiple charts.

**Emergent Properties:** A notable emergent property of our network is that the learned  $UV$  parametrization domains are consistent across different views of the same object (i.e., corresponding pixels in different views have similar  $UV$  coordinates) – and even across views of related objects in the same class. This is despite the  $UV$  domain maps only being indirectly supervised for consistency, through 3D reconstruction.

**Scope:** In this work, our focus is on *continuity, consistency, and 2D image–3D surface correspondences*. We focus on the case when the multi-view images have little overlap, a setting where traditional stereo matching techniques fail. Our method only requires supervision for the input views and their corresponding NOCS maps but does not require camera poses or ground truth  $UV$  parametrization. We note that the generated surfaces need not be watertight, and continuity at the seams between views is not guaranteed.

## 2 Related Work

There is a large body of work on object reconstruction which we categorize broadly based on the underlying shape representation.

**Voxels:** The earliest deep-learning-based methods predict a voxel representation of an object’s shape. Many of these methods are trained as generative models for 3D shapes, with a separate image encoder to obtain the latent code for a given image [14]. Later methods use more efficient data structures, such as octrees [40,45,36] to alleviate the space requirements of explicit voxels. Multiple views can also be aggregated into a voxel grid using a recurrent network [9]. Several methods use supervision in the form of 2D images from different view-points, rather than a 3D shape, to perform both single-view and multi-view voxel reconstruction [19,49,42,16]. These methods usually use some form of a

differentiable voxel renderer to obtain a 2D image that can be compared to the ground truth image. The quality gap of these methods to their counterparts that use 3D supervision is still quite large. Voxels only allow for a relatively coarse representation of a shape, even with the more efficient data representations. Additionally, voxels do not explicitly represent an object’s surface prompting the study of alternative representations.

**Point Clouds:** To recover the point cloud of an object instance from a single view, methods with 3D supervision [11,25] and without 3D supervision [17] have been proposed. These methods encode the input image into a latent code thus losing correspondences between the image pixels and the output points. Some methods establish a coarse correspondence implicitly by estimating the camera parameters, but this is typically inaccurate. A recent method reconstructs a point cloud of a shape from multiple views [6], but requires ground truth camera parameters. A large body of monocular or stereo depth estimation methods obtain a point cloud for the visible parts of the scene in an image, but do not attempt to recover the geometry of individual object instances in their local coordinate frames [3]. NOCS [43,39] obtains exact correspondences between 2D pixels and 3D points by predicting the 3D coordinates of each pixel in a canonical coordinate frame. NOCS can even be extended to reconstruct unseen parts of an object [39] (X-NOCS). All these approaches that output point clouds do not describe the connectivity of a surface, which has to be extracted separately – a classical and difficult geometry problem. We extend NOCS to directly recover continuous surfaces and consistently handle multiple views.

**Implicit Functions:** Poisson Surface Reconstruction [22,23] has long been the gold standard for recovering an implicit surface from a point cloud. More recently, data-driven methods have been proposed that model the implicit function with a small MLP [8,34,31], with the implicit function representing the occupancy probability or the distance to the surface. These methods can reconstruct an implicit function directly from a single image, but do not handle multiple views and do not establish a correspondence between pixels and the 3D space. PiFU [37] and DISN [48] are more recent methods that establish a correspondence between pixels and 3D space and use per-pixel features to parameterize an implicit function. Both single and multiple views can be handled, but the methods either require ground truth camera poses as input [37], or use a network to get a coarse approximation of the camera poses, giving only approximate correspondences [48]. Some recent works integrate the neural rendering with deep implicit functions [32,29], but they depend on the known camera information. Furthermore, to obtain an explicit surface from an implicit function, an expensive post-processing step is needed, such as Marching Cubes [30] or ray tracing.

**Parametric Surfaces or Templates:** Several methods attempt to directly reconstruct a parametric representation of a shape’s surface. These parametric representations range from class-specific templates [18,24], general structured templates [13], or more generic surface representations, such as meshes or continuous functions. Pixel2Mesh and its sequel [44,47] deform a genus-zero mesh based on local image features at each vertex, obtained by projecting the vertices to the

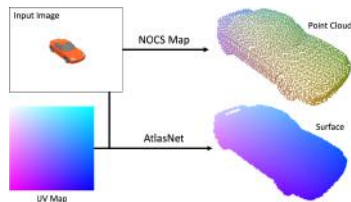
image plane(s). Camera parameters are assumed to be known for this projection. 3DN [46] deforms a given source mesh to approximate a single target image, using global features for both the source and the target, without establishing correspondences to the target pixels. Several methods use 2D images instead of 3D meshes as supervisory signal [21,28,35,20] using differentiable mesh renderers. This makes it easier to collect training data, but the accuracy of these methods still lags behind methods with 3D supervision. AtlasNet [15] represents shapes with continuous 2D patches that can be inferred from a single input image, or from a video clip [26]. Mesh DeformNet [33] introduces topology modification to AtlasNet. Similar to AtlasNet, we use a 2D patch as a *UV* parametrization, but we handle multiple non-adjacent views and establish correspondences between 2D pixels and 3D surface points.

### 3 Preliminaries

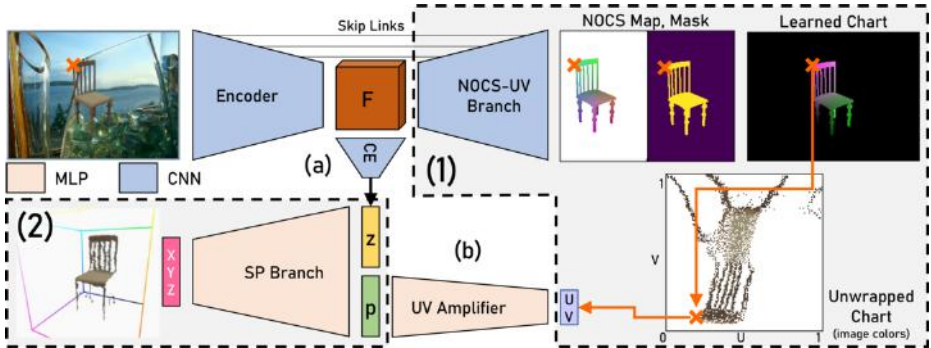
We build our approach upon two previous ideas that we describe below.

**(X-)NOCS:** Normalized object coordinate space (NOCS) is a canonicalized unit container space used for category-level reasoning of object pose, size, and shape [39,43]. Instances from a given object category are *normalized* for their position, orientation, and size, thus disentangling intra-category shape variation from the exact pose and size of instances. NOCS maps (see Fig. 2) are perspective projections of the 3D NOCS shape onto a specific camera and can be interpreted as **object-centered depth maps** that simultaneously encode mask and partial shape of the object. When used to predict 3D point cloud from images, NOCS maps retain correspondences from 2D pixels to 3D points, and can be used to transport image texture directly to 3D. X-NOCS is an extension of NOCS maps to also encode the occluded parts of a shape [39]. However, using NOCS maps for reconstruction results in a discontinuous point cloud.

**Surface Parametrization:** A two-manifold surface in 3D can be mapped to a 2D plane (*chart*) parametrized by two coordinates  $(u, v)$ . This *UV* parametrization of a 3D surface is widely used in computer graphics and, more recently, in 3D shape reconstruction [15,24]. The parameterization can be limited in expressing complex shapes, depending on the functional formulation used. For example, in typical CAD settings, low-degree polynomial or rational functions are used to represent the mappings. In our case, instead, we use a fully connected network to overcome the limitation of expressibility. A single map, however, still lacks the ability to describe complicated shapes with high-genus topology. Thus, multiple charts are often used, where multiple 2D planar patches are mapped by separate maps to a 3D surface – effectively partitioning the surfaces into parts, each of which is the image of a different map in the



**Fig. 2.** Given a single image, X-NOCS [39] reconstructs a point cloud preserving pixel–3D correspondence. AtlasNet [15] learns shape as a continuous surface.



**Fig. 3.** Single-View Single-chart Pix2Surf network architecture. The input image is processed using an encoder-decoder architecture to predict a NOCS map, object mask, and a learned chart (top right two channels, color coded). The Surface Parameterization branch takes sampled and amplified chart coordinates  $\mathbf{p}$  and a latent image code  $\mathbf{z}$  to predict the final continuous surface. Unwrapped chart (bottom right) refers to a visualization of foreground colors using the predicted two-channel learned chart (top right) as coordinate. The colors of the input image can be transported to all intermediate steps ( $\times$  and arrows).

chart. We show how a single chart can be used for 3D shape reconstruction while multiple charts allow consistent reconstruction over multiple views while still preserving pixel to 3D correspondence.

## 4 Pix2Surf: Pixels to Surface

Our goal is to predict a continuous and consistent parametric 3D surface for a novel object instance observed from one or more views. The word “continuous” parametric surfaces in our method refers to parametric  $C^0$  continuity (similar to AtlasNet [15]), i.e., any continuous trajectory over the chart space maps to a continuous curve in 3D space. Additionally, we would like to preserve correspondences between 2D pixels and 3D surface points. We first describe our approach for reconstructing a 3D surface from a single image using a single chart, and then generalize it to multiple views using an atlas.

### 4.1 Single-View Single-Chart Pix2Surf

At inference time, the single-view version of Pix2Surf takes an RGB image of an object observed from an arbitrary camera as input. We use a CNN to extract image features that compactly encode object shape. The features are then processed by two branches: (1) the **NOCS-UV branch** is a CNN that estimates a mask, a learned UV map, and a NOCS map and (2) the **Surface Parametrization (SP) branch** is an MLP that generates a continuous 3D surface. This single-view, single-chart architecture is shown in Fig. 3.

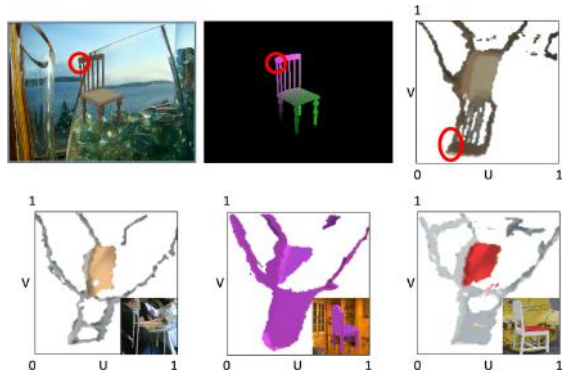
**(1) NOCS-UV Branch:** Similar to X-NOCS [39], we predict the NOCS map and mask that encode the partial shape of the object observed in the image. We use an encoder-decoder architecture building on top of SegNet [1] and VGG [38]. Our network uses skip connections and shares pool indices between the encoder and the decoder. The predicted NOCS maps and masks are the same size as the input image. During training, the object mask is supervised with a binary cross entropy loss and the NOCS map is supervised with an  $L^2$  loss. Note that the NOCS map here is *not* our final 3D output, but acts as an intermediate supervision signal for the network.

**Emergence of a Chart:** Different from previous work, we predict a 2-channel output in addition to the NOCS map and mask. These 2 channels are not explicitly supervised during training, so the network can predict any value between 0 and 1. However, when jointly trained with the other branches, we observe the **emergence** of a *learned chart* in these 2 channels (see Fig. 4). The network discovers how to unwrap an object shape onto a flat surface. **Remarkably, this learned chart is (almost) consistent across multiple views and even across instances.** During reconstruction, each image pixel’s learned chart coordinates are passed on to the SP branch. We show that using the learned chart coordinates is superior to using arbitrary UV coordinates like AtlasNet [15], or alternatively using the original image coordinates (Image2Surf, Sec. 5.1). Additionally, we preserve exact correspondences between input image pixels and the learned chart.

**(a) Code Extractor (CE):**

We use a small CNN to reduce the high dimensional feature map extracted by the encoder to make a more compact global code for the SP branch. This CNN contains two convolutional layers (512 and 1024 output channels), batch normalization, and ELU activation. The output is a latent code  $\mathbf{z}$  of size 1024 and is passed to the SP branch.

**(b) UV Amplifier:** Before we use the learned chart coordinates as an input to the SP branch, we process each UV coordinate with a *UV amplifier* MLP. The motivation for this comes from the information imbalance the two inputs to the SP branch – one input is the global latent code  $\mathbf{z}$  which



**Fig. 4.** Given an object image (row 1, col 1), our network predicts a 2-channel image without explicit supervision (row 1, col 2, color coded). Remarkably, the output of these two channels visualized in a  $UV$  space (row 1, col 3) show that the network has learned to unwrap the 3D shape onto a plane (corresponding patches shown in red). This unwrapping is consistent over multiple views, and even across multiple object instances (last row). For more unwrapped charts of cars and airplanes please see supplementary Fig. S.3.

has 1024 dimensions, while the UV coordinates would have only 2 dimensions. To overcome this, we *amplify* the UV coordinates to  $\mathbf{p}$  (256 dimensions) using a 3-layer MLP that progressively amplifies the 2 coordinates (2, 64, 128, 256). This allows the SP branch to make use of the image and UV information in a more balanced manner.

**(2) SP Branch:** Similar to AtlasNet [15], our surface parametrization (SP) branch takes the global latent code  $\mathbf{z}$  from the code extractor (CE) and the amplified coordinates  $\mathbf{p}$  as input and produces a continuous 3D position as the output. Note that the learned chart coordinates can be continuously sampled at inference time. The continuity of the output 3D surface emerges from our use of a continuous MLP mapping function between the uv coordinates and the output 3D positions [15]. Our SP branch is a MLP with 9 layers and skip connection every 2 layers (input: 1024+256, intermediate: 512, last: 3). Since we train on canonically oriented ShapeNet models, the predicted 3D surface also lies within the canonical NOCS container [43].

Our approach has three key differences to AtlasNet. First, we use a UV amplifier to transform the 2D  $UV$  coordinates to higher dimensions allowing better information balancing. Second, the learned chart is in direct correspondence with the pixels of the input image (see Fig. 4). This allows us to transport appearance information directly from the image to the 3D surface. Third, our sampling of the  $UV$  chart is learned by a network (NOCS-UV branch) instead of uniform sampling, which enables us to reconstruct complex topologies. Our inference processing allows us to sample any continuous point in the learned chart space within the predicted object mask allowing the generation of **continuous textured 3D surface**.

**Training:** The encoder and decoder CNNs are first initialized by training them on the NOCS map and mask prediction tasks using only the  $L^2$  and BCE losses. Subsequently, we jointly train the NOCS-UV and SP branches, code extractor, and UV amplifier end-to-end. The joint loss is given as,

$$\mathcal{L}_I = w_1 (w_n \mathcal{L}_n + w_m \mathcal{L}_m) + w_2 \mathcal{L}_s, \quad (1)$$

where  $\mathcal{L}_n$  and  $\mathcal{L}_m$  are the  $L^2$  NOCS map and BCE losses respectively,  $w_n, w_m$  are the weights for the NOCS map and mask prediction respectively, and  $w_1, w_2$  are the weights for the NOCS-UV and SP branches respectively. For the SP branch we supervise on  $K$  points sampled randomly from within the foreground mask. For each sampled point, a corresponding amplified chart coordinate  $\mathbf{p}$  is predicted without any supervision. This is concatenated with the global latent code  $\mathbf{z}$  to predict the final 3D surface position. Empirically, we found the best hyperparameters to be:  $K = 4096, w_1 = 0.1, w_2 = 0.9, w_n = 0.7, w_m = 0.3$ . The loss for the SP branch is given as,  $L_s = \frac{1}{K} \sum_{i=1}^K \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_2$ , where  $\mathbf{x}$  and  $\hat{\mathbf{x}}$  are the ground truth and predicted 3D surface position obtained from the 3D ShapeNet models (same as ground truth NOCS map values). During inference, we can predict a continuous 3D surface for any given image and its learned chart coordinate. Please see the supplementary document for more details on inference and final 3D model generation.

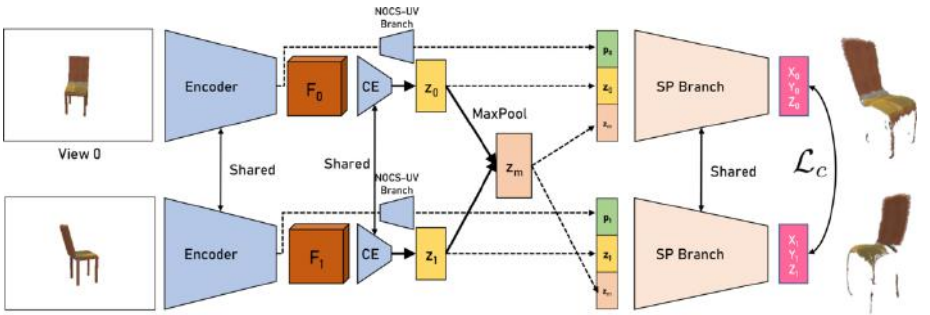


## 4.2 Multi-View Atlas Pix2Surf

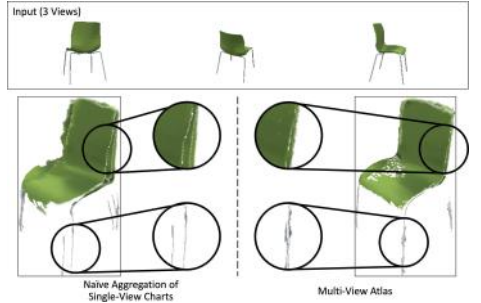
The method described above is suitable when we have a single view of the object. When multiple views are available, we could naively extend the single view network and combine the generated surfaces using a union operation. However, this leads to sharp discontinuities (Fig. 5). To overcome this issue, we propose a generalization of our single-view single-chart method to consistently aggregate 2D surface information from multiple views, using an *atlas* i.e., a separate learned chart (UV map) for each view. Fig. 6 shows an overview of our multi-view network. This design shares similarities with the single view network but has additional

multi-view consistency which is enforced both at the feature level through a feature pooling step, and using a consistency loss for better 3D surface generation.

**Multi-View Feature Pooling:** The goal of this step is to promote multi-view information sharing at the *feature level* (see Fig. 6). Different from the single view network, the latent codes  $\mathbf{z}_i$  extracted for each view  $i$  (using a shared encoder and code extractor) are maxpooled into a common shared multi-view latent code  $\mathbf{z}_m$ . Intuitively, this shared latent code captures the most salient information from each view.



**Fig. 6.** Multi-view atlas network architecture. The multi-view network allows multiple charts to be consistently aggregated. This network has two main features: (1) the MaxPool operation to pool features between views, and (2) a multi-view consistency loss  $\mathcal{L}_c$  that ensures corresponding points produce 3D surface points that are nearby. Only two views are shown in this figure, but we use multiple views during training and inference. The encoder, NOCS-UV branch, CE branch, and SP branches share weights.



**Fig. 5.** Given 3 views, naïve aggregation of individual charts leads to discontinuities or double surfaces (left). Our multi-view atlas method produces more consistent surfaces (right), for instance, at the legs and backrest.

**Atlas:** Similar to the single view network, we learn a chart for each view. The chart coordinates for each view  $\mathbf{p}_i$  are extracted using the NOCS-UV branch with weights shared between the views. Although the NOCS-UV branch weights are shared, one chart is predicted for each view – thus, we have an atlas. Note that the network is free to predict different chart coordinates for each view. However, we observe that similar parts of objects in different images map to similar locations on their respective charts (see Fig.S3 in supplementary document). This indicates that our network is **discovering the notion of image cross-view correspondence (note that this is different from 2D-3D correspondence)**. As in the single-view network, chart coordinates are passed through a shared UV amplifier.

We concatenate the shared latent code  $\mathbf{z}_m$  to each of the per-view latent codes  $\mathbf{z}_i$ . This concatenated multi-view code and the learned per-view chart coordinates  $\mathbf{p}_i$  are passed to the SP branch. The UV amplifier, code extractor and structure of the learned UV map are similar to the single view network.

**Multi-View Loss:** In addition to the  $L^2$  loss function on the 3D surface generated by the SP branch, we also have a multi-view consistency loss. This loss enforces corresponding points on multiple views to predict similar 3D surface positions. To obtain correspondence information at training time, we sample a random set of foreground points within the mask and find the exact match of the ground truth NOCS values of that pixel in the other input views. Note that this correspondence information is *not* provided as additional supervision – the ground truth NOCS maps already contain this information since corresponding points multiple views have the same NOCS position. Given these correspondences, the multi-view consistency loss for a pair of views is given as,  $\mathcal{L}_C = \frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \|\mathbf{x}_i - \mathbf{x}_j\|_2$ , where  $\mathbf{x}_{i,j}$  are the paired predicted xyz from two different views and the set  $\mathcal{P}$  contains all matched correspondence pair from these two views. During training, within each mini-batch, we sample multiple views per object and compute the loss for all possible pairs.

**Training:** The multi-view network is trained similar to the single view model. The NOCS-UV branch is first trained and subsequently the whole network is trained end-to-end. The loss function we use is  $\mathcal{L}_M = \mathcal{L}_I + \frac{w_3}{a} \sum_{j=0}^a \mathcal{L}_C$ , where  $a$  denotes the number of pairs of views within that batch, and  $w_3$  is the correspondence loss weight empirically set to 0.9. We set  $w_n, w_m$  to 0.1 inside  $\mathcal{L}_I$ . Please see the supplementary document for more details on inference and final 3D model generation.

## 5 Experiments

We present extensive experimental comparison of Pix2Surf with several recent single- and multi-view reconstruction methods, and validate our design choices. We do so by focusing on the 3C properties (consistency, correspondence and continuity) for visible surface reconstruction (Sec. 5.1). Since we learn a strong prior over shapes, we can also estimate surfaces that are hidden in the input image (Sec. 5.2). For the details of training, inference, and evaluation metrics,

and also ablations, more comparisons, and results with real images, refer to the supplementary document.

**Dataset:** For quantitative comparisons, we use ShapeNetPlain [39] dataset which consists of 5 random views for each shape in ShapeNet [4] with a white background. For additional robustness to the background found in real-world images, we train Pix2Surf on ShapeNetCOCO [39] which consists of 20 random views of each ShapeNet object with a random background from MS COCO [27]. We use this dataset for all qualitative results and for real-world results. Each shape category is trained separately in all experiments.

**Experimental Setting:** We follow the experimental setup of X-NOCS [39]. Our ground truth for each input image is the point cloud represented by the NOCS map (or X-NOCS map for hidden surface) provided in the dataset [39]. The outputs of all methods are converted to a NOCS map (using the ground truth camera pose) allowing us to compute metrics even for partial shapes. Multi-view experiments use all 5 views in the dataset to reconstruct a surface, using the same dataset as the single-view experiments. All metrics are computed per-view and then averaged up, making the single- and multi-view values comparable in our quantitative experiments.

**Metrics:** We quantify the quality of reconstructed surfaces with several metrics. The **reconstruction error** of predictions is computed as the Chamfer distance [2,12] between the estimated surface and the ground truth NOCS map (interpreted as a point cloud). To obtain points on a reconstructed surface, we convert it into a NOCS map using the ground truth camera pose.

In addition to the accuracy of reconstructed surfaces, we quantify the 3C properties of a surface with the following metrics. The 2D–3D **correspondence error** measures the accuracy of the estimated correspondence between input pixels and 3D points on the reconstructed surface. The error for each foreground pixel is the distance between the estimated 3D location of the pixel and the ground truth location. Unlike the Chamfer distance, this uses the 2D–3D correspondence to compare points. We average over all foreground pixels to obtain the correspondence error of a surface. The **multi-view consistency error** was defined in Sec. 4.2 as the 3D distance between corresponding points in different views. We average the distance for a given point over all pairs of views that contain the point. Corresponding points are found based on the ground truth NOCS map. We measure **discontinuity** based on the surface connectivity. While the continuity of Pix2Surf is induced by our use of a continuous MLP as mapping from uv space to 3D space [15], the mapping from the *input image* to the 3D space should *not* be  $C^0$ -continuous everywhere, due to self occlusions and boundaries of the 3D shape. The reconstructed surface should have the same  $C^0$  discontinuities as the ground truth surface. We define a  $C^0$  discontinuity as large difference in the 3D locations of the neighboring pixels in a NOCS map (above a threshold of 0.05). We take a statistical approach to measure the surface connectivity, by computing a histogram over the 3D distances between neighboring pixels that are discontinuous. The discontinuity score is the correlation of this histogram to a histogram of the ground truth surface. A higher score indicates a distribution

**Table 1.** Visible surface reconstruction. We compare our method to a baseline and three state-of-the-art methods evaluating reconstruction accuracy and the 3C properties. The top half of the table shows single-view reconstruction, the bottom half is multi-view reconstruction. Note how Pix2Surf is close to the top performance in each of the metrics, while all other methods have significant shortcomings. The **Recons. Error** and **Correspond. Error**, **Consistency Error** are all multiplied by  $10^3$ .

	<b>Recons. Error</b> ↓				<b>Correspond. Error</b> ↓				<b>Consistency Error</b> ↓				<b>Disconti. Score</b> ↑			
	car	chair	plane	avg.	car	chair	plane	avg.	car	chair	plane	avg.	car	chair	plane	avg.
Im.2Surf	2.23	3.81	2.66	2.90	<b>8.49</b>	9.54	8.76	8.93	13.08	12.55	10.75	12.13	0.46	0.39	0.35	0.40
X-NOCS	2.25	2.95	2.08	2.43	12.82	8.63	8.93	10.13	18.93	12.00	10.59	13.84	0.59	<b>0.47</b>	0.59	0.55
AtlasNet	<b>1.54</b>	3.36	3.15	2.68	—	—	—	—	—	—	—	—	0.68	0.39	0.64	0.57
<b>Pix2Surf</b>	<b>1.67</b>	<b>1.91</b>	<b>1.61</b>	<b>1.73</b>	9.52	<b>5.79</b>	<b>7.19</b>	<b>7.50</b>	<b>12.72</b>	<b>7.75</b>	<b>8.48</b>	<b>9.65</b>	<b>0.69</b>	0.43	<b>0.65</b>	<b>0.59</b>
X-NOCS	2.89	2.80	2.19	2.63	14.30	9.48	8.95	10.91	22.18	14.26	11.65	16.03	<b>0.67</b>	<b>0.48</b>	0.54	0.56
P2M++	2.88	5.59	3.24	3.90	—	—	—	—	—	—	—	—	0.67	0.36	0.63	0.55
<b>Pix2Surf</b>	<b>1.41</b>	<b>1.78</b>	<b>1.38</b>	<b>1.52</b>	<b>8.49</b>	<b>5.84</b>	<b>7.06</b>	<b>7.13</b>	<b>10.98</b>	<b>6.65</b>	<b>7.50</b>	<b>8.38</b>	0.66	0.43	<b>0.66</b>	<b>0.58</b>

of discontinuities that is more similar to the ground truth surface. Note that continuity is a property induced from method design itself, and the score can penalize the over-smooth case from methods that produces continuous prediction.

## 5.1 Visible Surface Reconstruction

We compare the quality of single- and multi-view reconstructions to one baseline [Image2Surf (single-view)], and three state-of-the-art methods [AtlasNet [15] (single-view), X-NOCS [39] (single- and multi-view), Pixel2Mesh++ [47] (multi-view)]. Note that Pix2Surf deals with a more challenging problem compared to AtlasNet and Pixel2Mesh++: (1) we predict 2D–3D correspondences (AtlasNet does not), and (2) we do not require camera geometry information as input (Pixel2Mesh++ does). In this section, we only focus on reconstructing *visible* surfaces, but we also report hidden surface generation in the next section.

The single-view performance of each method in all of our metrics is shown in the first four rows of Table 1, and the multi-view performance in the last three rows. Metrics are comparable across single- and multi-view methods. For each of the four metrics, we show the performance on each dataset category, and an average over all categories.

**Image2Surf:** This baseline is similar to Pix2Surf, but takes image  $UV$  coordinates (normalized by predicted mask) as input to the  $UV$  amplifier instead of the learned  $UV$  chart, i.e., the input image is the chart. We observe that it is hard for the network to learn depth discontinuities, resulting in over-smoothed occlusion boundaries (see supplementary document). The over-smoothing is reflected in a high reconstruction error, and particularly low discontinuity correlation score. This comparison justifies our design to include a learned  $UV$  chart.

**X-NOCS:** This is a state-of-the-art reconstruction method that predicts a 3D point cloud, i.e., a 3D point for each foreground pixel. Since X-NOCS has no notion of surface connectivity, there is no coordination between neighboring points, resulting in poor reconstruction accuracy and noisy output point clouds (see Fig. 7).

Note that the output point cloud from X-NOCS can capture the right discontinuity. However, it can only produce discrete noisy point cloud instead of continuous surfaces.

**AtlasNet:** This method also uses an explicit surface parametrization, giving it a low reconstruction error on the Car category. However, since the parametrization is not learned and has

a fixed layout and connectivity, the reconstruction error increases significantly for categories with more complex shapes and topologies, such as Chair and Airplane. Correspondence and multi-view consistency are not evaluated, since AtlasNet lacks pixel-to-point correspondences and works only for a single view.

**Pixel2Mesh++:** This method deforms a given starting mesh in a coarse-to-fine approach to approximate an object shown from multiple views. In each refinement step, a mesh vertex is deformed based on a small image neighborhood around the projection of the vertex in each view. Unlike in our method, ground truth camera positions need to be known for this projection. The fixed connectivity and topology of the starting mesh results in a higher reconstruction error. Since correspondence and multi-view consistency are trivial given a ground truth camera model, we do not evaluate these properties.

Unlike the previous methods, **Pix2Surf** learns a continuous parametrization of the surface that does not have a fixed topology or connectivity. This gives us more flexibility to approximate complex surfaces, for instance, to correctly place holes that can model  $C^0$  discontinuities. This explains our high discontinuity correlation scores which also benefits the accuracy of reconstruction and 2D-3D correspondence. In the multi-view setting, Pix2Surf shares information across the views, improving the overall reconstruction accuracy. For example, surfaces that are only visible at a very oblique angle in one view can benefit from additional views. Our use of a consistency loss additionally ensures an improvement of the multi-view consistency over the baselines, and a lower consistency error compared to single view Pix2Surf (Fig. 5). We observe that Pix2Surf is the only method that has top performance on all quality metrics (reconstruction and 3C properties), all other methods reconstruct surfaces that fall short in at least some of the metrics.



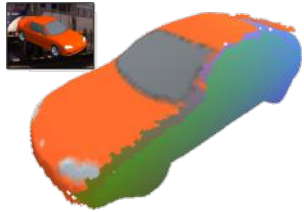
**Fig. 7.** Our results (left) compared with surface-agnostic X-NOCS (right), visualized with image connectivity. Pix2Surf produces significantly smoother results.

## 5.2 Hidden Surface Generation

Since Pix2Surf learns a strong prior of the shapes it was trained on, we can generate plausible estimates for surfaces in parts of the object that are not directly visible in the image (see Fig. 8). Similar to X-NOCS, we represent a 3D object with two layers: a visible layer that we reconstruct in the experiments described previously, and a hidden layer denoting the last intersection of camera rays [39]. Pix2Surf can be easily extended to reconstruct hidden surface farthest from the camera by adding additional output channels to the NOCS-UV branch. The rest of the architecture remains the same with the learned  $UV$  parametrization

**Table 2.** We compare the reconstruction error of visible and hidden surfaces (trained jointly) for Pix2Surf and X-NOCS [single view (sv.) and multi-view (mv.)]. The learned parametrization of Pix2Surf also benefits from hidden surface generation, and the additional reconstruction of the hidden surface does not adversely affect the accuracy of the visible surfaces.

	Visible Error ↓				Hidden Error ↓			
	car	chair	plane	avg.	car	chair	plane	avg.
X-NOCS (sv.)	2.25	2.95	2.08	2.43	1.86	3.34	2.25	2.48
X-NOCS (mv.)	2.89	2.80	2.19	2.63	3.11	3.32	2.03	2.82
<b>Pix2Surf</b>	<b>1.66</b>	<b>2.01</b>	<b>1.66</b>	<b>1.78</b>	<b>1.52</b>	<b>2.47</b>	<b>1.77</b>	<b>1.92</b>



**Fig. 8.** Pix2Surf can reconstruct both visible (textured) and hidden parts (color coded).

additionally also learning about the hidden surface. In Table 2, we show our performance when jointly reconstructing the visible and hidden surfaces from an image. We compare to both the single- and multi-view version of X-NOCS on all categories. The improvement in accuracy for our method shows that hidden surfaces benefits from our learned parametrization as well. Comparing the performance of the visible surface reconstruction to Table 1, we see that the joint reconstruction of visible and hidden surfaces does not significantly decrease the reconstruction accuracy of the visible surfaces.

## 6 Conclusion

We have presented Pix2Surf, a method for predicting 3D surface from a single- or multi-view images. Compared with the previous work, Pix2Surf simultaneously achieves three properties in the prediction: **continuity** of the surface, **consistency** across views, and pixel-level **correspondences** from the images to the 3D shape. By attaining these properties, our method enables the generation of high-quality parametric surfaces, readily integrating the output surfaces from multi-views, and lifting texture information from images to the 3D shape. In future work, we will explore ways of guaranteeing continuity even across different views and improving the quality of mapped textures. Another interesting direction is to exploit the intermediate learned chart as a container for material properties. A longer-term goal would be to investigate how the network can generalize across multiple categories.

**Acknowledgements:** We thank the anonymous reviewers for their comments and suggestions. This work was supported by a Vannevar Bush Faculty Fellowship, NSF grant IIS-1763268, grants from the Stanford GRO Program, the SAIL-Toyota Center for AI Research, AWS Machine Learning Awards Program, UCL AI Center, and a gift from the Adobe.

## References

1. Badrinarayanan, V., Kendall, A., Cipolla, R.: Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* (2017)
2. Barrow, H.G., Tenenbaum, J.M., Bolles, R.C., Wolf, H.C.: Parametric correspondence and chamfer matching: Two new techniques for image matching. In: *Proceedings of the 5th International Joint Conference on Artificial Intelligence - Volume 2*. pp. 659–663. IJCAI’77, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA (1977), <http://dl.acm.org/citation.cfm?id=1622943.1622971>
3. Bhoi, A.: Monocular depth estimation: A survey. *arXiv preprint arXiv:1901.09402* (2019)
4. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012* (2015)
5. Chen, R., Han, S., Xu, J., Su, H.: Point-based multi-view stereo network. In: *Proc. of ICCV* (2019)
6. Chen, R., Han, S., Xu, J., Su, H.: Point-based multi-view stereo network. In: *Proc. of ICCV* (2019)
7. Chen, Z., Zhang, H.: Learning implicit fields for generative shape modeling. In: *Proc. of CVPR* (2019)
8. Chen, Z., Zhang, H.: Learning implicit fields for generative shape modeling. In: *Proc. of CVPR* (2019)
9. Choy, C.B., Xu, D., Gwak, J., Chen, K., Savarese, S.: 3D-R2N2: A unified approach for single and multi-view 3d object reconstruction. In: *Proc. of ECCV* (2016)
10. Deprelle, T., Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: Learning elementary structures for 3d shape generation and matching. In: *Proc. of NeurIPS* (2019)
11. Fan, H., Su, H., Guibas, L.: A point set generation network for 3d object reconstruction from a single image. In: *Proc. of CVPR* (2017)
12. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 605–613 (2017)
13. Genova, K., Cole, F., Vlastic, D., Sarna, A., Freeman, W.T., Funkhouser, T.: Learning shape templates with structured implicit functions. In: *Proc. of ICCV* (2019)
14. Girdhar, R., Fouhey, D., Rodriguez, M., Gupta, A.: Learning a predictable and generative vector representation for objects. In: *Proc. of ECCV* (2016)
15. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: Atlasnet: A papier-mch approach to learning 3d surface generation. In: *Proc. of CVPR* (2018)
16. Henzler, P., Mitra, N., Ritschel, T.: Escaping plato’s cave using adversarial training: 3d shape from unstructured 2d image collections. In: *Proc. of ICCV* (2019)
17. Insafutdinov, E., Dosovitskiy, A.: Unsupervised learning of shape and pose with differentiable point clouds. In: *Proc. of NeurIPS* (2018)
18. Kanazawa, A., Tulsiani, S., Efros, A.A., Malik, J.: Learning category-specific mesh reconstruction from image collections. In: *Proc. of ECCV* (2018)
19. Kar, A., Häne, C., Malik, J.: Learning a multi-view stereo machine. In: *Proc. of NeurIPS* (2017)
20. Kato, H., Harada, T.: Learning view priors for single-view 3d reconstruction. In: *Proc. of CVPR* (2019)

21. Kato, H., Ushiku, Y., Harada, T.: Neural 3d mesh renderer. In: Proc. of CVPR (2018)
22. Kazhdan, M., Bolitho, M., Hoppe, H.: Poisson surface reconstruction. In: Proc. of the Eurographics symposium on Geometry processing (2006)
23. Kazhdan, M., Hoppe, H.: Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)* (2013)
24. Kulkarni, N., Gupta, A., Tulsiani, S.: Canonical surface mapping via geometric cycle consistency. In: Proc. of ICCV (2019)
25. Lin, C.H., Kong, C., Lucey, S.: Learning efficient point cloud generation for dense 3d object reconstruction. In: Proc. of AAAI (2018)
26. Lin, C.H., Wang, O., Russell, B.C., Shechtman, E., Kim, V.G., Fisher, M., Lucey, S.: Photometric mesh optimization for video-aligned 3d object reconstruction. In: Proc. of CVPR (2019)
27. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: Proc. of ECCV (2014)
28. Liu, S., Li, T., Chen, W., Li, H.: Soft rasterizer: A differentiable renderer for image-based 3d reasoning (2019)
29. Liu, S., Saito, S., Chen, W., Li, H.: Learning to infer implicit surfaces without 3d supervision. In: *Advances in Neural Information Processing Systems*. pp. 8295–8306 (2019)
30. Lorensen, W.E., Cline, H.E.: Marching cubes: A high resolution 3d surface construction algorithm. In: *SIGGRAPH* (1987)
31. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: Proc. of CVPR (2019)
32. Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 3504–3515 (2020)
33. Pan, J., Han, X., Chen, W., Tang, J., Jia, K.: Deep mesh reconstruction from single rgb images via topology modification networks. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 9964–9973 (2019)
34. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: DeepSDF: Learning continuous signed distance functions for shape representation. In: Proc. of CVPR (2019)
35. Petersen, F., Bermano, A.H., Deussen, O., Cohen-Or, D.: Pix2vex: Image-to-geometry reconstruction using a smooth differentiable renderer. *arXiv preprint arXiv:1903.11149* (2019)
36. Richter, S.R., Roth, S.: Matryoshka networks: Predicting 3D geometry via nested shape layers. In: Proc. of CVPR (2018)
37. Saito, S., Huang, Z., Natsume, R., Morishima, S., Kanazawa, A., Li, H.: Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In: Proc. of ICCV (2019)
38. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
39. Sridhar, S., Rempe, D., Valentin, J., Bouaziz, S., Guibas, L.J.: Multiview aggregation for learning category-specific shape reconstruction. In: Proc. of NeurIPS (2019)
40. Tatarchenko, M., Dosovitskiy, A., Brox, T.: Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In: Proc. of ICCV (2017)
41. Tatarchenko, M., Richter, S.R., Ranftl, R., Li, Z., Koltun, V., Brox, T.: What do single-view 3D reconstruction networks learn? In: Proc. of CVPR (2019)



42. Tulsiani, S., Zhou, T., Efros, A.A., Malik, J.: Multi-view supervision for single-view reconstruction via differentiable ray consistency. In: Proc. of CVPR (2017)
43. Wang, H., Sridhar, S., Huang, J., Valentin, J., Song, S., Guibas, L.J.: Normalized object coordinate space for category-level 6d object pose and size estimation. In: Proc. of CVPR (2019)
44. Wang, N., Zhang, Y., Li, Z., Fu, Y., Liu, W., Jiang, Y.G.: Pixel2mesh: Generating 3d mesh models from single rgb images. In: Proc. of ECCV (2018)
45. Wang, P.S., Liu, Y., Guo, Y.X., Sun, C.Y., Tong, X.: Adaptive O-CNN: A Patch-based Deep Representation of 3D Shapes. In: SIGGRAPH Asia (2018)
46. Wang, W., Ceylan, D., Mech, R., Neumann, U.: 3dn: 3d deformation network. In: Proc. of CVPR (2019)
47. Wen, C., Zhang, Y., Li, Z., Fu, Y.: Pixel2mesh++: Multi-view 3d mesh generation via deformation. In: Proc. of ICCV (2019)
48. Xu, Q., Wang, W., Ceylan, D., Mech, R., Neumann, U.: DISN: Deep implicit surface network for high-quality single-view 3D reconstruction. In: Proc. of NeurIPS (2019)
49. Yan, X., Yang, J., Yumer, E., Guo, Y., Lee, H.: Perspective transformer nets: Learning single-view 3d object reconstruction without 3d supervision. In: Proc. of NeurIPS (2016)

## Supplementary Material

### S.1 Overview

In this supplementary, we provide additional details about our training (Sec. S.2) and inference setups (Sec. S.3), and details of our evaluation metrics (Sec. S.4). We provide an extended qualitative comparison of our method to the Image2Surf baseline (Sec. S.5), ablations (Sec. S.6) and for visible surface generation on real-world data (Sec. S.7). We show additional qualitative results for hidden surface generation (Sec. S.9) and also provide more visual results for Pix2Surf (Sec. S.8) and more qualitative comparison to Pixel2Mesh++ [47] and AtlasNet [15] (Sec. S.10).

### S.2 Training Details

For the **Single-View** case, we train our network in two phases. In the first phase, we train the NOCS-UV branch with a learning rate of  $1e-4$ , using the NOCS Map and the object mask as supervision. In the second phase, we add the remaining SP branch and train end-to-end until convergence, with a learning rate of  $1e-4$  for cars and  $3e-5$  for planes and chairs, and using the losses described in Sec. 4.1.

For the **Multi-View** case, we have found that pre-training with the single-view architecture, before switching to the full multi-view architecture results in better initialization. For this purpose, we start by passing the feature  $z_m$ , directly to the SP branch without max-pooling multiple views. After pre-training, we switch to the multi-view architecture as described in Sec. 4.2, by max-pooling the  $z_m$  features of all views, and concatenating both this max-pooled multi-view feature, and the single-view feature  $z_m$  for the current view as input to the MLP. To better fuse multi-view information for learned chart prediction, the feature map in the middle of CNN encoder and decoder also follows above fusion operation. We randomly pick 5 views as input during multi-view training. For our multi-view consistency loss, we need to identify corresponding pixels in different views. We sample pixels in each view as in the single-view case and find corresponding pixels based on their distance in NOCS coordinates. Two pixels are in correspondence if their NOCS distance is less than  $1e-3$ .

We separately train on each object category of our dataset.

### S.3 Inference Details

One significant advantage of our explicit **continuous** parametric surface prediction is that we can sample the results at any resolution (e.g. points or vertices). We generate our final predictions at a regular grid of samples in the unwrapped uv chart, obtaining a 3D location for each sample (obtained from the SP-Branch). Since we have exact correspondence to pixels of the input image, each sample also has a color value (or interpolated color value in super-resolution). Samples corresponding to background pixels are masked out. To create a mesh, we can connect neighboring foreground samples with edges. All visual results of our method in the paper are generated using this approach. We provide more details.

*Identifying foreground regions in the unwrapped chart.* Unlike AtlasNet, the shape and topology of the unwrapped surface in our chart is learned by the NOCS-UV branch, which gives the reconstructed surface more flexibility to represent arbitrary shapes and topologies. To identify foreground regions in the uv space of the unwrapped chart, we map the the learned image-space foreground mask to uv space. Directly unwrap the mask by learned-uv map (two channel output from NOCS-UV branch) results in pixel cloud with holes in uv space. To solve this issue, we up-sample the image-space mask and learned-uv map from its original resolution of  $240 \times 320$  by a factor of 4 using linear interpolation before mapping mask to uv space. To avoid interpolating across  $C^0$  discontinuities of the surface, we only interpolate neighboring pixels that are mapped to similar uv locations (i.e., the gradient of their uv coordinates is below a threshold). We then map the up-sampled mask to uv space (resolution of  $128 \times 128$ ) by the up-sampled learned-uv map. Finally we up-sample the mask in uv space to the desired resolution (in paper we use  $512 \times 512$ ).

In uv space, we additionally post-process the unwrapped foreground mask by closing small holes using morphological operations. Finally, we remove outliers using the predicted 3D locations (quarried from SP-Branch) of each mask sample. A sample of the foreground mask is classified as outlier if the distance in 3D space to its nearest neighbor is larger than a threshold  $t$ . In practice, we use  $t = 0.03$  for chairs and  $t = 0.02$  for cars and airplanes. Similar outlier removing operation is also applied to image-space mask before identifying foreground regions.

*Texturing the unwrapped chart.* Similar to the mask, directly unwrapping the image-space color values to the uv space results in a sparse set of irregular color samples in uv space. We can interpolate these samples to obtain the color value at any point in uv space by interpolating the  $k$  nearest neighbors (we use  $k = 4$  for our results).

## S.4 Evaluation Metrics

We now define the evaluation metrics used in the paper.

**A common surface representation:** Before evaluating our metrics, we convert the results of all methods to a common format to avoid biasing our results due to different surface representations. We convert all output representations to the NOCS-Map format defined in X-NOCS [39] using the ground truth camera model. The NOCS map  $\mathcal{P}$  samples the reconstructed surface from a single viewpoint, giving a point cloud where each sample has a 2D pixel coordinate  $p$  and a 3D location  $x$ . The 3D location is defined in a canonical coordinate frame that is shared across views and across instances of the same shape category. For multi-view reconstructions, we create one NOCS-Map for each viewpoint, compute the metrics on each NOCS-Map, and average the results over all views. As AtlasNet [15] ground truth is not in the same ShapeNet version as ShapeNet-Plain [39], we first scale the AtlasNet results to have the same bounding box diagonal as the ground truth 2-intersection X-NOCS maps point cloud, and then align the lower left corner of the bounding box.

The **Reconstruction Error** is measured as the 2-Way-Chamfer-Distance between the ground truth NOCS-Map  $\mathcal{P}_1$  and predicted NOCS-Map  $\mathcal{P}_2$ :

$$E_{\text{rec}} = \frac{1}{|\mathcal{P}_1|} \sum_{x_i \in \mathcal{P}_1} \min_{y_j \in \mathcal{P}_2} \|x_i - y_j\|_2^2 + \frac{1}{|\mathcal{P}_2|} \sum_{y_j \in \mathcal{P}_2} \min_{x_i \in \mathcal{P}_1} \|x_i - y_j\|_2^2.$$

The reconstruction error for hidden surfaces in Table 2 of paper is computed in the same way, but using NOCS-Maps of the hidden surfaces.

The **Correspondence Error** is measured as the squared distance between the predicted 3D location  $x_i$  and the ground truth location  $y_i$  of the same pixel:

$$E_{\text{corr}} = \frac{1}{|\mathcal{M}|} \sum_{p_i \in \mathcal{M}} \|x_i - y_i\|_2^2.$$

We only evaluate pixels  $p_i \in \mathcal{M}$  that are both in the predicted and ground truth foreground masks.

**Consistency Error** is based on the squared distance between the predicted 3D locations of corresponding pixels in different views. For each pair of views  $a$  and  $b$ , we identify corresponding pairs of pixels  $(p_i^a, p_j^b)$  as pairs having a similar ground truth 3D location in NOCS:  $\|y_i^a - y_j^b\|_2 < \epsilon$ . In practice, we set  $\epsilon = 0.001$ . We then average the squared distance between the predicted 3D locations  $x_i^a$  and  $x_j^b$  of all corresponding pixel pairs  $\mathcal{P}_{\text{corr}}^2$ :

$$E_{\text{cons}} = \frac{1}{|\mathcal{P}_{\text{corr}}^2|} \sum_{(p_i^a, p_j^b) \in \mathcal{P}_{\text{corr}}^2} \|x_i^a - x_j^b\|_2^2.$$

With the **Discontinuity Score**, we take a statistical approach to measure the correctness of the surface connectivity. While the continuity of implicit or parametric surface is a property induced by representation and method design, we need to make sure the continuity is correct, i.e. no over-smooth results. We compute statistics of the  $C^0$  discontinuities in the predicted surface, and measure the similarity to the same statistics computed on the ground truth surface. The statistics are based on the 3D distance  $\|x_i - x_j\|_2$  of neighboring foreground pixels  $p_i$  and  $p_j$ . Pixels with a large difference are likely to lie on the border of a  $C^0$  discontinuity of the predicted surface. We compute a histogram  $h$  of this 3D distance over all neighboring pixels:

$$h_i = |\{(p_i, p_j) \in \mathcal{P}_{\text{neighbors}}^2 \mid t_i \leq \|x_i - x_j\|_2 < t_{i+1}\}|,$$

where  $t_i$  are the boundaries of the histogram bins and  $\mathcal{P}_{\text{neighbors}}^2$  is the set of all neighboring pixel pairs. We use a 4-connected neighborhood and choose 20 bins with bin edges spaced uniformly in  $[0.05, \sqrt{3}]$ . We measure the similarity of two histograms as the correlation of their normalized bins:

$$S_{\text{cont}} = \frac{1}{\sum_k h_k \sum_j h_j^{\text{gt}}} \sum_i h_i h_i^{\text{gt}}$$

Note that unlike the other errors we use as evaluation metrics, this is defined as a score, where higher values imply more accurate discontinuity of the reconstructed surface.

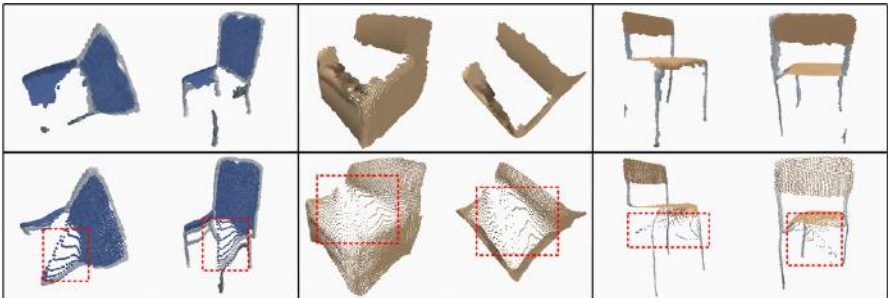
## S.5 Qualitative Comparison to Image2Surf

We show more qualitative comparisons between our baseline Image2Surf and Pix2Surf in paper Sec. 5.1. Image2Surf has a fatal problem to make “cut” around the occlusion boundary (i.e., wrong  $C^0$  discontinuities), which is reflected both in the red rectangle in Fig. S1 and discontinuity score in Table 1 in paper.

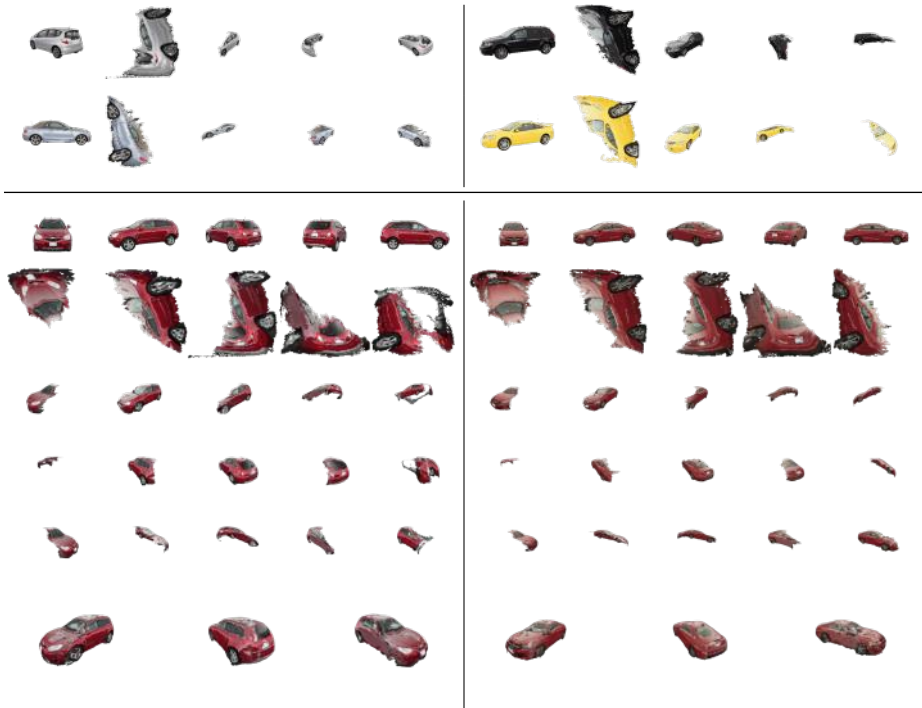
## S.6 Ablations

We provide an analysis and justification of several key design choices. First, we analyze the importance of using a learned  $UV$  chart instead of a fixed chart like Image2Surf. As seen in Sec. 5.1 and Table 1, Pix2Surf outperforms Image2Surf on all categories for reconstruction error. Second, we analyze the utility of multi-view feature pooling and consistency loss. As seen in Table 1 (rows 4–7), these two features significantly improve performance. We also justify the use of intermediate NOCS map regression by the NOCS-UV branch, and the need for the UV amplifier (Table S1). We do so by examining networks without these two components. For the NOCS map ablation, we train the network from scratch without pretraining the NOCS-UV branch, and for the UV amplifier ablation, we directly input the learned UV coordinates to the SP branch and increase the dimension of a latent image code to 256. When conducting the experiments on the chair category, the results (Table S1) show that these components help learn better reconstructions. **Table S1.** We experimentally verify the usefulness of NOCS map regression and the UV amplifier. NOCS map regression provides intermediate supervision while the UV amplifier balances information. Here we report average reconstruction error computed on the *visible* part (equal training epochs for all methods).

	No UV Amp.	No NOCS	Pix2Surf
Chair	10.37	3.64	<b>2.61</b>



**Fig. S1.** Qualitative Comparison to Image2Surf. The first row are the results of Pix2Surf and the second row are for Image2Surf. Each instance is viewed from 2 different viewpoints. Image2Surf wrongly connects disjoint parts and results in strong distortions, which are solved by Pix2Surf’s learned chart.



**Fig. S2.** Real world image generalization. The top part is single-view visualization: input image, unwrapped chart with texture and 3 viewpoints of the reconstruction for each instance. The bottom part is multi-view aggregation visualization. For every instance, each row is: input images, unwrapped charts with texture, 3 viewpoints for each view’s result separately and finally multi view aggregation.

## S.7 Qualitative Results on Real-World Data

We show more results for generalization to real world data mentioned in paper. In Fig. S2, we show single- and multi-view results for Pix2Surf that is trained on ShapeNet COCO and inference on real world car. Note that the texture in each view separately is better than the multi-view aggregation. This is caused by the different light condition from different viewpoints. As our main concern in this paper is not to fuse the texture from multiple views, we leave the improvement of the texture to future works.

## S.8 More Results

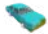

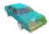



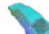










































































Figure S3 shows more results of Pix2Surf including the learned UV map (as shown in Figure 4) and reconstruction outputs of both single-view and multi-view architectures. See the caption for the details.



**Fig. S3.** Single-view and multi-view Pix2Surf reconstruction results. The results for each object are presented in three rows. The first row shows five input views. Note that we do not have camera parameters for any of these views. The second row shows the per-view UV space that is generated by the multi-view variant of Pix2Surf. The UV space is not directly constrained by any loss; the flattening of the objects that we can observe and the large degree of consistency between different views is an emergent property of our network. In the third row, we show, from left to right, (a) the reconstructed 3D surface obtained by merging Pix2Surf single-view reconstructions (SV), (b) the Pix2Surf multi-view reconstruction (MV), and (c) the ground truth reconstruction (GT). The last three columns show the same results from a different viewpoint. Note the reduction in the number of gaps and surface discontinuities when comparing the multi-view to the single-view results.

## S.9 Qualitative Results for Hidden Surface Generation



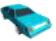



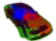







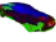







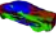







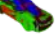

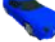

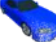



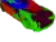







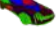







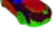













































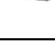











The following table provides more visual results of Pix2Surf Two-Intersection version (Sec 5.2 in paper), and comparison with X-NOCS [39]. Pix2Surf can easily be extended to capture the invisible surface and is more accurate and smooth than X-NOCS.

View 1			View 2			View 3		
Pix2Surf (sv)	X-NOCS (sv)	Ground Truth	Pix2Surf (sv)	X-NOCS (sv)	Ground Truth	Pix2Surf (sv)	X-NOCS (sv)	Ground Truth
								
								
								
								
								
								
								
								
								



## S.10 Qualitative Comparisons

The following table demonstrates qualitative comparisons among our Pix2Surf (both single-view and multi-view architectures), AtlasNet[15], and Pixel2Mesh++ [47]. The colors in AtlasNet results show different output patches.

View 1			View 2			View 1	
Single View	Multi-View	Ground Truth	Single View	Multi-View	Ground Truth	Atlas Net	Pixel2 Mesh++
							
							
							
							
							
							
							
							
							
							
							
							
							
							

View 1			View 2			View 1	
Single View	Multi-View	Ground Truth	Single View	Multi-View	Ground Truth	Atlas Net	Pixel2 Mesh++
