

Learning Shape Templates with Structured Implicit Functions

Kyle Genova^{1,2} Forrester Cole² Daniel Vlasic² Aaron Sarna² William T. Freeman² Thomas Funkhouser^{1,2}

¹Princeton University ²Google Research

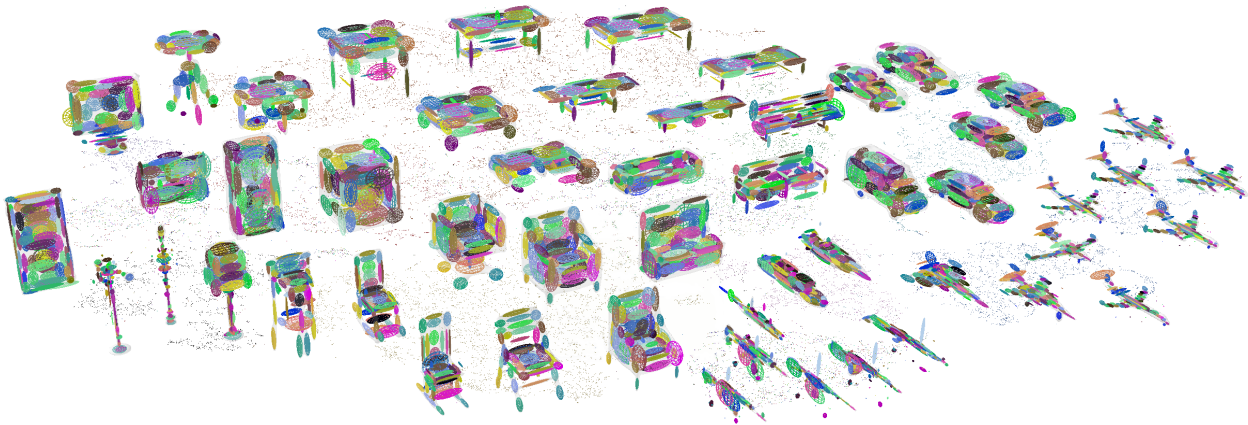


Figure 1. Shapes from the ShapeNet [8] database, fit to a structured implicit template, and arranged by template parameters using t-SNE [52]. Similar shape classes, such as airplanes, cars, and chairs, naturally cluster by template parameters.¹

Abstract

Template 3D shapes are useful for many tasks in graphics and vision, including fitting observation data, analyzing shape collections, and transferring shape attributes. Because of the variety of geometry and topology of real-world shapes, previous methods generally use a library of hand-made templates. In this paper, we investigate learning a general shape template from data. To allow for widely varying geometry and topology, we choose an implicit surface representation based on composition of local shape elements. While long known to computer graphics, this representation has not yet been explored in the context of machine learning for vision. We show that structured implicit functions are suitable for learning and allow a network to smoothly and simultaneously fit multiple classes of shapes. The learned shape template supports applications such as shape exploration, correspondence, abstraction, interpolation, and semantic segmentation from an RGB image.

1. Introduction

Fitting a 3D shape template to observations is one of the oldest and most durable vision techniques [43]. Templates offer a concise representation of complex shapes and

a strong prior for fitting. They can be used to directly correspond and compare shapes, and supervised learning approaches may be applied to correspond the template and a photograph [54, 5]. In order to fit a wide range of shapes, however, multiple, hand-made templates are usually required, along with a procedure for choosing the appropriate one [13].

The goal of this paper is to construct a general shape template that fits any shape, and to learn the parameters of this template from data. We view a shape as a level set of a volumetric function and approximate that function by a collection of shape elements with local influence, a formulation we term a *structured implicit* function. The template itself is defined by the number of and formula for the shape elements, and the template parameters are simply the concatenation of the parameters of each element. An example of this type of representation is the classic *metaballs* method [3], but more sophisticated versions have been proposed since [57, 4, 38].

Given a template definition, we show that a network can be trained to fit the template to shapes with widely varying geometry and topology (Figure 1). Critically, the net-

¹See templates.cs.princeton.edu for video, supplemental, and high resolution results.

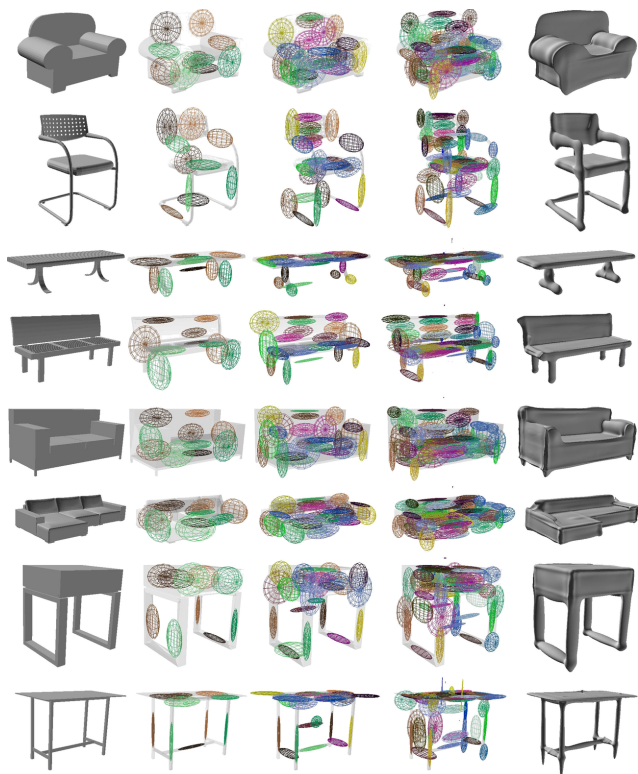


Figure 2. Templates fit to a variety of geometry and topology. Middle columns: three shape templates trained across classes with 10, 25, and 100 elements, respectively. Right: surface reconstruction of the implicit function defined by the 100 element template. Note how the structure of each template is consistent between shapes.

work learns a fitting function that is smooth: the template parameters of similar shapes are similar, and vary gradually through shape-space (Figure 2). Further, we show that the network learns to associate each shape element with similar structures in each shape: for example, the tail fin of an aircraft may be represented by one element, while the left wingtip may be represented by another. This consistency allows us to interpolate shapes, estimate vertex correspondences, or predict the influence region of a given element in a 2D image, providing semantic segmentation of shapes.

The closest related work to ours is the volumetric primitive approach of Tulsiani, et al. [51]. Like that work, we aim to learn a consistent shape representation with a small number of primitives. We expand on their work by specifying the surface as a structured implicit function, rather than as a collection of explicit surface primitives. This change allows for an order of magnitude increase in the number of shape elements, allowing our template to capture fine details.

Our method is entirely self-supervised and requires only a collection of shapes and a desired number of shape elements (N , usually 100). The output template is concise ($7N$ values) and can be rendered or converted to a mesh using techniques such as raytracing or marching cubes [32].

2. Related Work

There is a long history of work on shape analysis aimed at extracting templates or abstract structural representations for classes of shapes [19, 18, 26, 35, 58].

Primitive Fitting: Fitting of basic primitives is perhaps the oldest topic in 3D computer vision, beginning with Roberts [43] and continuing to today [2, 20, 31, 29]. These methods focus on explaining individual observations with primitives, and do not necessarily provide consistency across different input shapes, so they cannot be used for the correspondence, transfer, and exploration applications targeted in this paper.

Part Segmentation: Others have studied how to decompose mesh collections into consistent sets of semantic parts, either through geometric [14] or learned methods [1, 12, 21, 27]. These methods differ from ours in that they depend on labeled examples to learn the shapes and arrangements of *semantic parts* within specific classes. In contrast, we aim to learn a *structural* template shape for any class without human input.

Template Fitting: The most related techniques to ours are methods that explicitly fit templates to shapes [7]. The templates can be provided by a person [13, 39], derived from part segmentations [60, 28, 12], or learned automatically [24, 51, 59, 60]. Previous work generally assumes an initial set of primitives or part structure is given prior to learning. For example, Kim et al. [24] proposed an optimization to fit an initial set of box-shaped primitives to a class of 3D shapes and used them for correspondence and segmentation. Part structure is assumed by [60].

Others have learned shape templates with a neural network. In Zou et al. [61], a supervised RNN is trained to generate sets of primitives matching those produced by a heuristic fitting optimization. Sharma et al. [46] use reinforcement learning to decompose input shapes into a CSG parse tree. Like our approach, this approach does not require additional training data, but CSG trees are unsuitable for many template applications.

Tulsiani et al. [51] proposed a neural network that learned placements for a small number (3 to 6) of box primitives from image or shape inputs, without additional supervision. Our method builds on this approach, but greatly expands the number and detail of the shape elements, allowing for the precise shape associations required for correspondence and semantic segmentation applications.

Implicit Shape Representations: Decades ago, researchers in computer graphics proposed representing shapes with sets of local shape functions [42, 3]. The most common form is a summation of polynomial or Gaussian basis functions centered at arbitrary 3D positions, sometimes called *metaballs* [3], blobby models [36], or soft objects [57].

Property	Voxel	Octree	Point	Mesh	Deep	Ours
Interpret	+	+	+	+	-	+
Concise	-	+	+	+	-	+
Surface	+	+	-	+	-	+
Volume	+	+	-	-	+	+
Topology	+	+	-	-	+	+
Deform	-	-	+	+	-	+

Table 1. Comparison of desirable properties of various 3D representations, rated as suitable (+) or unsuitable (-). From top to bottom: is the representation interpretable to humans; concise in storage; capable of representing surfaces and volumes; allows topological changes; and supports smooth deformation. Structured implicit functions are suitable in all properties. “Deep” refers to methods that represent a volumetric function as a deep neural network [40, 47].

Other forms include convolution surfaces [4] and partition of unity implicits [38]. These representations support compact storage, efficient interior queries, arbitrary topology, and smooth blends between related shapes, properties that are particularly useful for our application of predicting template shapes.

Shape Representations for Learning: Recently, several deep network architectures have appeared that encode observations (color images, depth images, 3D shapes, etc.) into a latent vector space and decode latent vectors to 3D shapes. Our work follows this approach. We argue that our structured implicit representation is superior for template learning compared to decoding voxels [6, 55, 56], sparse-voxel octrees [50], points [11], meshes [15, 22, 53], box primitives [51], signed-distance function estimators [40], or indicator function estimators [34].

Table 1 compares the properties of these representations. Compared to points, implicit surfaces are superior because they provide a clearly-defined surface. Compared to meshes, implicit surfaces can continuously adapt to arbitrary topology. Structured implicit functions are most similar to voxel grids since both implicitly represent a surface. Unlike voxel grids, they provide a sparse representation of shape, though octree techniques can provide sparse representations of voxels. The major difference for our work is that our shape elements can be moved and transformed in a smooth way to, for example, track gradual changes in airplane wing shape across a shape collection. By contrast, two similar, but slightly transformed shapes will have entirely different voxel representations.

Techniques have recently been proposed to directly approximate volumetric functions such as signed-distance fields or indicator functions using deep neural networks [40, 34, 47]. Compared to these approaches, structured implicit functions are light weight, easily interpretable, and provide template geometry that can be modified or transformed by later processing.

3. Structured Implicit Shape Representation

We assume each input shape can be modeled as a watertight surface bounding an interior volume (real-world meshes usually must be processed to satisfy this assumption, see Sec. 4.2). We aim to represent this surface as the ℓ level set of a function $F(\mathbf{x}, \Theta)$, where \mathbf{x} is a 3D position and Θ is a vector of template parameters. In the structured implicit formulation, F is the sum of the contributions of a fixed number of shape elements with local influence, labeled $i \in [N]$, where N is their count. Each element is a function f_i defined by its parameter vector θ_i (making Θ simply the concatenation of θ_i):

$$F(\mathbf{x}, \Theta) = \sum_{i \in [N]} f_i(\mathbf{x}, \theta_i) \quad (1)$$

The specific version of shape elements we adopt are *scaled axis-aligned anisotropic 3D Gaussians*. Here, θ_i consists of a scale constant c_i , a geometric center $\mathbf{p}_i \in \mathcal{R}^3$, and per-axis radii $\mathbf{r}_i \in \mathcal{R}^3$.

$$f_i(\mathbf{x}, \theta_i) = c_i \exp \left(\sum_{d \in \{x, y, z\}} \frac{-(\mathbf{p}_{i,d} - \mathbf{x}_d)^2}{2\mathbf{r}_{i,d}^2} \right) \quad (2)$$

Intuitively, one can think this of representation as a set of squished or stretched 3D blobs. We found this set of parameters to be the minimum necessary to achieve good results. More sophisticated shape elements, such as full multivariate Gaussians, or even windowed quadric functions [38], would likely improve results, but we do not experiment with those here.

Because all constants c_i are negative, we have that $f_i(\mathbf{x}, \theta_i) < 0$ and thus $F(\mathbf{x}, \Theta) < 0, \forall \mathbf{x} \in \mathcal{R}^3$. Therefore we pick a negative isolevel ℓ and define the surface S to be its crossing:

$$S = \{\mathbf{x} \in \mathcal{R}^3 : F(\mathbf{x}, \Theta) = \ell\} \quad (3)$$

We set $\ell := -0.07$, which was chosen by grid search. The reason that the constants are negative rather than positive is to maintain the convention that function values inside the surface should be less than ℓ , while values outside the surface should be greater than ℓ . This leads to a convenient binary outside/inside test for points x :

$$F(\mathbf{x}, \Theta) > \ell \quad (4)$$

For most experiments presented here, we use $N = 100$. Because each shape element has seven parameters, the total dimensionality of our representation is a fixed $7N = 700$ floating point values.

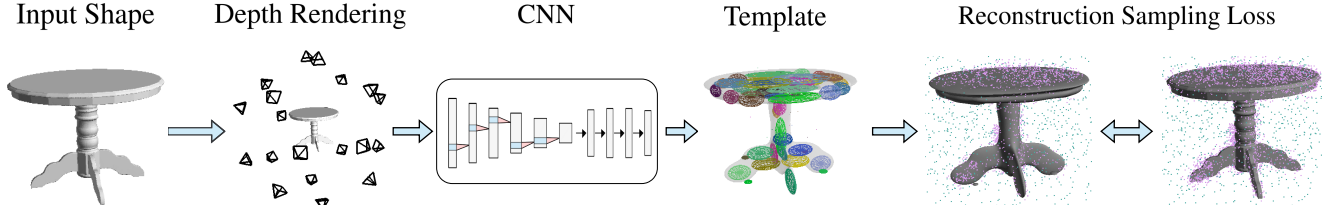


Figure 3. An overview of our method. The input to our system is a mesh. We render a stack of depth images around the mesh, and provide these as input to an early-fusion CNN. The output of the CNN is a vector with fixed dimensionality. This vector is interpreted as a shape template with parameters that define an implicit surface. Next, we sample points near the ground truth surface and also uniformly in space. A classification loss enforces that each sample point is correctly labeled as inside/outside by the surface reconstruction.

4. Template Learning

We propose a learning framework (Figure 3) to train a neural network to fit the shape template to data. The network’s goal is to find the template parameters Θ that best fit a 3D shape, where the loss penalizes the amount of predicted shape that is on the wrong side of the ground truth inside/outside border. We render multiple depth images of the mesh from fixed views to provide 3D input to the network. Our network has a feed-forward CNN architecture and predicts the entire parameter vector Θ at once with a fully connected layer. During training, we choose sparse sample locations in 3D and evaluate our loss function at those locations with a classification loss. The details of this procedure are described in the rest of this section.

Note that although fitting consistency is vital to our applications, we do not directly enforce similar shapes to have similar template parameters; the network arrives at a smooth fitting function without intervention. We hypothesize that, as a matter of optimization, the smooth solution is “easier” for the network to learn, but analyzing the causes of this behavior is an engaging direction for future work.

4.1. Architecture

In order to learn the template, we first need to encode the input 3D shape. There are a variety of network architectures for encoding 3D shape; options include point networks [41], voxel encoders [33], or multi-view networks [49]. Because voxel encoders can be computationally expensive, and point cloud encoders discard surface information, we opt for a multi-view encoding network. Our network takes a stack of 20 depth images rendered from the vertices of a dodecahedron as input, as in [23]. The network contains 5 convolutional layers followed by 4 fully connected layers.

The final fully connected layer is linear and maps to the template parameter vector Θ , which in our experiments is usually 700-D. Even though we use an encoder/decoder style architecture, there is no heavy decoding stage: the code vector is our explicit representation. We experimented with alternative “decoding” architectures, such as an LSTM that predicts each shape element in succession. We found

the LSTM architecture to perform better in some cases, but it took much longer to train, and was not able to scale easily to large numbers of shape elements.

4.2. Data Preparation

Before training, we must preprocess the input meshes to make them watertight. This step is important primarily because our loss function requires a ground truth inside/outside classification label.

In order to do the watertight conversion, we first convert the meshes to a 300^3 sparse voxel representation [37]. We flood fill the octree to determine inside/outside, then extract the isocontour of the volume to produce the watertight mesh. We generate 100,000 random samples uniformly in the bounding box of the mesh, and compute 0/1 inside/outside labels. We additionally compute 100,000 samples evenly distributed on the surface of the mesh.

We also render depth maps of the watertight meshes. For each mesh, we render 20 depth images at uniformly sampled viewing directions as input to the network. The (depth maps, labeled samples) pairs are the only data used for learning.

4.3. Loss

The goal of our loss function is only to measure deviation from the input shape; we assume that our representation will naturally create a smooth template due to its structure. In order to accurately reconstruct the surface, we employ three individual loss functions, described in detail in the following sections. L_U and L_S are classification losses ensuring that the volume around the ground truth shape is correctly classified as inside/outside. These losses were inspired by recent work on implicit function learning [34, 9]. L_C enforces that all of the shape elements contribute to the reconstruction. The total loss function is a weighted combination of the three losses:

$$L = w_U L_U + w_S L_S + L_C \quad (5)$$

L_C has no weight here because it contains two subclasses with different weights w_a and w_b .

As our losses compare the structured implicit value $F(\mathbf{x}, \Theta)$ to indicator function labels (0 inside, 1 outside), we formulate a soft classification boundary function to better facilitate gradient learning:

$$G(\mathbf{x}, \Theta) = \text{Sigmoid}(\alpha(F(\mathbf{x}, \Theta) - \ell)) \quad (6)$$

where α controls the sharpness of the boundary, and is set to 100 as determined by grid search.

4.3.1 Uniform Sample Loss L_U

If $F(\mathbf{x}, \Theta)$ correctly classifies every point in the volume according to the ground truth shape boundary, then it has perfectly reconstructed the ground truth. To measure the classification accuracy, we choose (x, y, z) coordinates uniformly at random in the bounding box of the ground truth mesh. We evaluate $F(\mathbf{x}, \Theta)$ at these locations, and apply a loss between the softened classification decision G , and the ground truth class label, which is 0 inside and 1 outside:

$$L_U(\mathbf{x}, \Theta) = \begin{cases} \beta G(\mathbf{x}, \Theta)^2 & \mathbf{x} \text{ inside} \\ (1 - G(\mathbf{x}, \Theta))^2 & \mathbf{x} \text{ outside} \end{cases} \quad (7)$$

At each training batch we randomly select 3,000 of the precomputed 100,000 points to evaluate the loss. β accounts for the inside/outside sample count differences.

4.3.2 Near Surface Sample Loss L_S

While the uniform sample loss is effective, it is problematic because it prioritizes surface reconstruction based on the fraction of the volume that is correct. The network can easily achieve 99%+ correct volume samples and still not visually match the observation. In particular, thin structures are unimportant to a volumetric loss but subjectively important to the reconstruction. To improve performance, we sample proportionally to surface area, not volume. We additionally want to ensure that the network is not biased to produce an offset surface, so the loss should be applied with similar weight on both the positive and negative side of the surface boundary.

In order to achieve these goals, we implemented the following algorithm. For each of the 100,000 surface samples, a ray is cast in each of the positive and negative normal directions away from the surface point. Because the mesh is watertight, at least one of the two samples must intersect the surface. The minimum of these two intersection distances is chosen, and truncated to some threshold. We sample a point along either normal direction with probability inversely proportional to the squared distance from the surface and proportional to the minimum intersection distance. The output samples roughly satisfy both of our goals: no thin structures are missed, regardless of their volume, and there is an equal sampling density on both sides of the surface.

This loss function, L_S , is identical to L_U (see Equation 7) except for the sample locations where it is applied. Note that L_S and L_U are not redundant with one another. Because L_S only contains samples very near the surface, it does not on its own enforce that the network keep free space clear of spurious shapes. We found it most effective to use a weighted combination of both losses, using L_S to do hard example mining, and L_U to ensure that free space around the shape remains clear.

4.3.3 Shape Element Center Losses L_C

One problem with the loss so far is that it is only concerned with the final composite function $F(\mathbf{x}, \Theta)$. If shape elements do not affect F , they also don't affect the loss. This "death" of shape elements can easily happen over time, since elements are randomly initialized and some are likely to be far from the ground truth surface. Their contribution to L_U and L_S is small, and there is no incentive for the network to use them. Our solution to this problem is to apply a third loss L_C , the center classification loss. This loss enforces that all predicted centers must lie on the inside of the predicted shape and within the ground truth bounding box:

$$L_C(\mathbf{x}, \Theta) = \begin{cases} w_a G(\mathbf{x}, \Theta)^2 & \mathbf{x} \in B \\ w_b \sum_d \max(0, B_L - \mathbf{x}_d, B_U - \mathbf{x}_d)^2 & \mathbf{x} \notin B \end{cases} \quad (8)$$

Above, w_a and w_b are hyperparameters balancing the two cases, which are in different units. B is the axis aligned bounding box of the ground truth shape, which has a lower coordinate B_L and an upper coordinate B_U . It states that if the predicted center \mathbf{x} is inside the ground truth bounding volume (where L_U will be applied, keeping free space empty), then \mathbf{x} must also be inside the predicted surface. On the other hand, if \mathbf{x} is outside the ground truth bounding box, then it should be directly encouraged to move inside the bounding volume because it can't be useful to the template from that distance.

5. Experiments

We conduct experiments to demonstrate important properties of the shape template: it accurately fits a wide variety of shapes, fits similar shapes with similar templates, can be used to find 3D-to-3D and 2D-to-3D correspondences, and can be fit from RGB images alone. We train and test on ShapeNet Core V2 [8], using the dataset split defined by 3D-R2N2 [10]. We show results trained on both the full dataset (Sections 5.1, 5.3, 5.4) and trained per-class (Sections 5.2 and 5.5). Identical hyperparameters we used to train all templates.

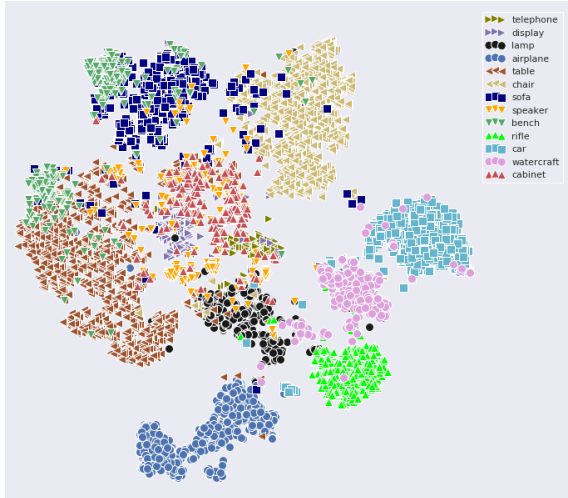


Figure 4. t-SNE visualization of template parameters on ShapeNet test set, colored by shape class labels. Note the clean clustering of most classes. Mixed clusters are also intuitive, e.g. mixing between tables, benches, and sofas.

5.1. Clustering by Template Parameters

A desirable property of a template fitting procedure is that similar shapes are fit with similar template parameters. Figure 4 shows a t-SNE [52] visualization of the template parameter vectors Θ for the ShapeNet test set, colored by ShapeNet class labels. Several classes of shapes (airplanes, rifles, cars) are neatly clustered by their template parameters. Other classes are mixed, but in intuitive ways: some benches look like tables, other benches look like sofas, and some sofas look like chairs. Cabinets, speakers, and displays are all essentially boxes, so they have similar template parameters.

5.2. Comparison to Volumetric Primitives

The closest alternative approach to ours is the volumetric primitives of Tulsiani, et al. [51]. We provide a detailed comparison between our template shapes and their shape abstractions using results generously provided by the authors. For this comparison we trained one fitting network per shape class, not one network for all classes, to match the procedure of [51]. Figure 5 shows representative results for examples from the ShapeNet training set, with 10, 25, and 100 shape elements (see supplemental material for the full set of results). In comparison to volumetric primitives (Figure 5 a), our templates (b-d) are more detailed, have higher consistency, and better reflect the structure of the input mesh (f).

5.3. Single-View RGB Prediction and Labeling

Figure 6 shows qualitative results demonstrating predictions from photographs of ShapeNet-style objects. To predict the template parameters from an RGB image, we apply

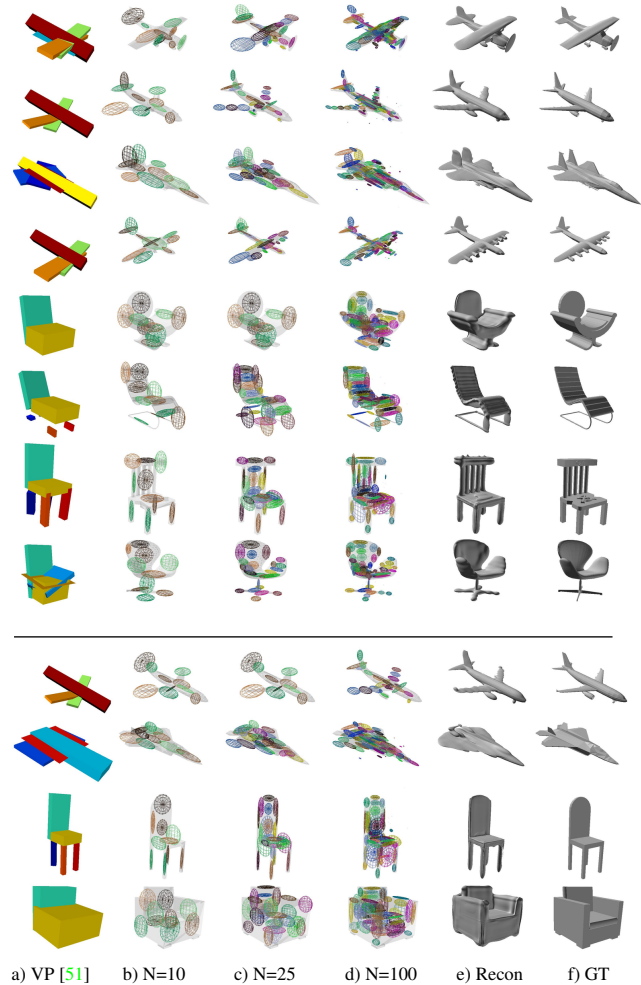


Figure 5. Comparison to Volumetric Primitives [51]: (a) volumetric primitives result; (b-d) templates computed with our method for 10, 25, and 100 elements; (e) surface reconstruction from the template in (d); (f) ground truth surface mesh. Shapes above the line come from our training set, while the shapes below the line are from our test set.

a similar technique to CNN purification [30] or network distillation [17] and train a second network that regresses from RGB to the template parameters already found through our 3D-to-3D training scheme. The training data for this network is synthetic OpenGL renderings of the ShapeNet training set, with camera angles chosen randomly from a band around the equator of the shape.

Because the template is consistent, we can go further than overall 3D shape prediction and predict correspondence between pixels in the image and the influence regions of individual shape elements (Figure 6, right). Each element tends to produce a particular part of each shape: the i^{th} element might produce the tail fin of an airplane, while the j^{th} might produce the wingtip. Because of this consistency,

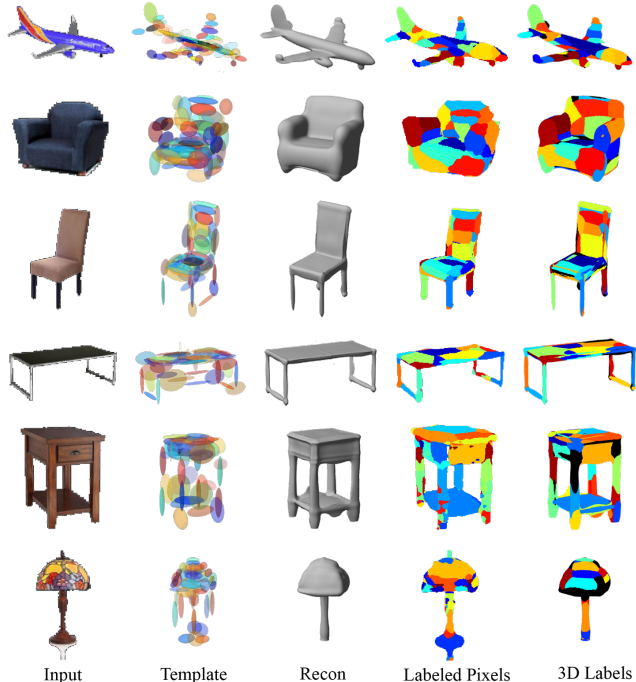


Figure 6. Template fitting and labeling from photographs. From left to right: input image with background removed, fit template, corresponding isosurface, image pixels labeled by the highest-value shape element, corresponding 3D regions labeled by the highest-value element. Regions in 3D not found by the image labeling network are black. The labeling performs well for easily oriented shapes (top rows), and worse for shapes with rotational symmetries (bottom rows). Note that the labeling is based entirely on the template, without additional region or part labels.

a semantic segmentation network [44] can be trained to label pixels by the index of the shape element with maximum weight at that pixel. The result is a segmentation of the image into 3D regions, without additional region or part labels. One limitation of this approach is that the template learning does not take into account object symmetry. Shapes with natural orientations, such as airplanes and chairs, are successful, while shapes without fronts and backs, such as the lamp and nightstand, confuse the network.

Similar techniques have been used for human body pose prediction [54, 5] using hand-made templates, but to our knowledge, we are the first to use a learned template.

5.4. Shape Correspondence

The learned template is consistent across shapes of the same class, meaning that the same elements will influence equivalent shape parts (e.g. airplane wings). This property can be exploited to find correspondences between different shapes. We present one automatic approach to achieve that. First, we use our network to compute the template configuration Θ of each shape we want to correspond. Then,

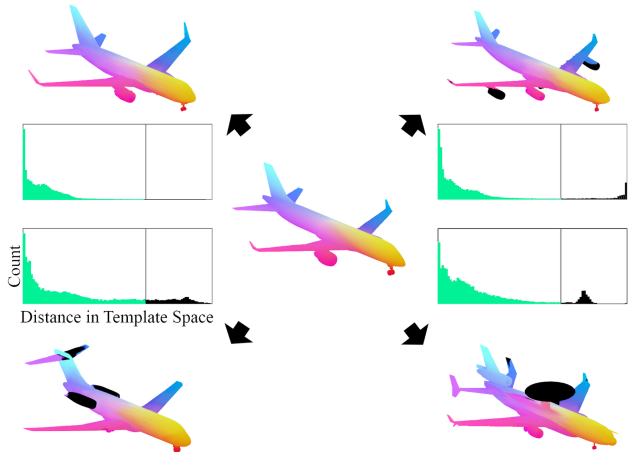


Figure 7. Transferring per-vertex colors from source airplane (center) to target airplanes (corners). Vertices are corresponded to their nearest neighbor in template space. Matching colors indicate corresponding vertices, while black regions have no corresponding vertices in the source. The histograms plot the proportion of nearest-neighbor distances that produce good matches (green) and outliers (black, distance > 0.65). Outliers include extra wing and tail engines, landing gear, and a radar dome, all missing on the source airplane. Correspondences were computed for resampled ShapeNet meshes from the training set of the multi-class network.

for each vertex \mathbf{v} , we compute its template coordinates. The template coordinates consist of three numbers for each shape element. Those are computed by subtracting the element’s center from the vertex position, dividing each coordinate by the corresponding element radius (improving correspondence between elongated and squashed elements), then scaling that vector to be of length $F(\mathbf{v}, \Theta)$. The direction of each per-element vector helps geometrically localize the vertex, while its length denotes the influence of that element. Finally, the cosine distance between template coordinates can be used to find the closest target vertex for each source vertex, as visualized in Figure 7.

5.5. Shape Interpolation

Another benefit of the structured template is the ability to geometrically interpolate between shapes. One can simply take multiple templates, blend all their parameters with per-template weights, and obtain an in-between template and hence the in-between shape. Figure 8 shows interpolation between four sofas computed this way.

5.6. RGB Single View 3D Reconstruction

While exact shape reconstruction is not the focus of our work, we compared the reconstruction accuracy of the template surface with the output of 3D-R2N2 [10], Point Set Generation Network [11], and Pix2Mesh [53]. The inputs are single RGB images of unknown camera orientation, so we use the distillation approach from Section 5.3. The

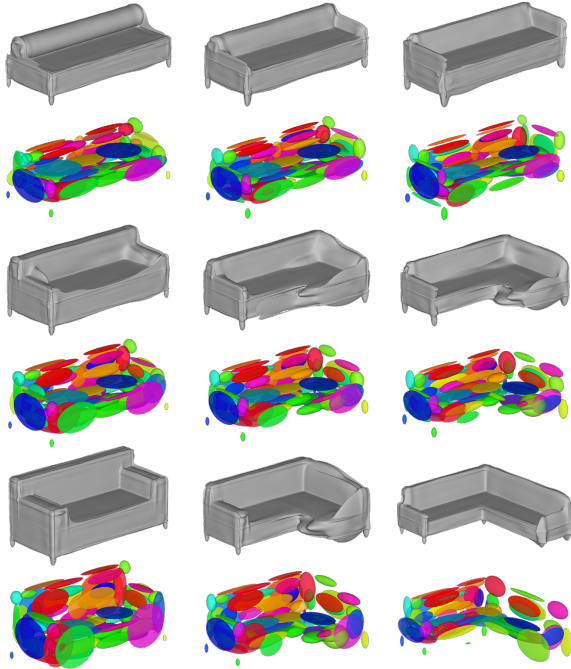


Figure 8. Linear interpolation between shape templates and corresponding surfaces of four sofas. The reconstructed shapes are at the four corners, pairwise blends are in-between the neighboring corners, and the average of all four shapes is in the center. The blends were obtained from training data of a single-class network.

Threshold	τ				2τ			
	R2N2	PSG	P2M	Our	R2N2	PSG	P2M	Our
plane	41	68	71	69	63	81	81	86
bench	34	49	58	62	49	69	72	82
cabinet	50	40	60	40	65	67	77	64
car	38	51	68	47	55	78	84	70
chair	40	42	54	40	55	64	70	64
monitor	34	40	51	42	48	64	67	65
lamp	32	41	48	32	44	59	62	52
speaker	45	32	49	29	58	57	66	50
firearm	28	70	73	72	47	83	83	88
sofa	40	37	52	42	53	63	70	70
table	44	53	66	40	59	73	79	61
cellphone	42	56	70	56	61	80	83	79
watercraft	37	51	55	49	52	71	70	75
mean	39	49	60	48	55	70	74	70

Table 2. F-score (%) on ShapeNet test set, with $\tau = 10^{-4}$ as in [53]. Higher numbers are better.

train/test split is from 3D-R2N2. Our shape representation has only 700 degrees of freedom, compared with $32^3 = 32768$ DoF for the 3D-R2N2 voxel grid, $1024 * 3 = 3072$ DoF for points generated by PSG, and $2466 * 3 = 7398$ DoF for the vertices of the Pix2Mesh mesh. Despite having many fewer degrees of freedom, the template surface



a) Template fit b) Reconstruction c) Input mesh

Figure 9. Shapes with angled parts, sharp creases, and thin structures are difficult for our method to learn.

reconstruction accuracy is similar to competing approaches (Table 2).

5.7. Limitations

Our method has several limitations apparent in Figure 9, which exhibits several failure cases. First, since our representation comprises of a small number of axis-aligned functions, it has limited ability to represent detailed, sharp, or angled structures (e.g., creases or corners). Second, since it learns to classify sides of a surface boundary, it struggles to reconstruct razor thin structures. Finally, since it uses a fixed number of shape elements (e.g., 100), it does not produce a template with 1-to-1 mapping to semantic shape components. We believe these limitations could be addressed with alternative (higher-order, non axis-aligned) local functions, distance-based loss functions, supervised training, and/or network architecture search.

6. Conclusion

This paper investigates using structured implicit functions to learn a template for a diverse collection of 3D shapes. We find that an encoder-decoder network trained to generate shape elements learns a template that maps detailed surface geometry consistently across related shapes in a collection with large shape variations. Applications for the learned template include shape clustering, exploration, abstraction, correspondence, interpolation, and image segmentation. Topics for future work include learning to generate higher-order and/or learned shape elements, deriving semantically meaningful shape elements via supervised learning, and using structured implicit functions for other applications such as 3D reconstruction.

7. Acknowledgements

We acknowledge ShapeNet [8], 3D-R2N2 [10], and the Stanford Online Products Dataset [48] for providing training and evaluation data for our method. We also thank the authors of Volumetric Primitives [51] for providing extended results from their method for our comparisons. We thank Avneesh Sud for helpful discussions and comments.

References

- [1] N. Araslanov, S. Koo, J. Gall, and S. Behnke. Efficient single-view 3d co-segmentation using shape similarity and spatial part relations. In *German Conference on Pattern Recognition*, pages 297–308. Springer, 2016. 2
- [2] I. Biederman. Recognition-by-components: a theory of human image understanding. *Psychological review*, 94(2):115, 1987. 2
- [3] J. F. Blinn. A generalization of algebraic surface drawing. *ACM transactions on graphics (TOG)*, 1(3):235–256, 1982. 1, 2
- [4] J. Bloomenthal and K. Shoemake. Convolution surfaces. *ACM SIGGRAPH Computer Graphics*, 25(4):251–256, 1991. 1, 3
- [5] F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *Computer Vision – ECCV 2016*, Lecture Notes in Computer Science. Springer International Publishing, Oct. 2016. 1, 7
- [6] A. Brock, T. Lim, J. M. Ritchie, and N. Weston. Generative and discriminative voxel modeling with convolutional neural networks. *arXiv preprint arXiv:1608.04236*, 2016. 3
- [7] R. Brunelli. *Template matching techniques in computer vision: theory and practice*. John Wiley & Sons, 2009. 2
- [8] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 1, 5, 8, 11
- [9] Z. Chen and H. Zhang. Learning implicit fields for generative shape modeling. *CoRR*, abs/1812.02822, 2018. 4
- [10] C. B. Choy, D. Xu, J. Gwak, K. Chen, and S. Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. 5, 7, 8, 11
- [11] H. Fan, H. Su, and L. J. Guibas. 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*, pages 605–613, 2017. 3, 7
- [12] N. Fish, M. Averkiou, O. Van Kaick, O. Sorkine-Hornung, D. Cohen-Or, and N. J. Mitra. Meta-representation of shape families. *ACM Transactions on Graphics (TOG)*, 33(4):34, 2014. 2
- [13] V. Ganapathi-Subramanian, O. Diamanti, S. Pirk, C. Tang, M. Nießner, and L. Guibas. Parsing geometry using structure-aware shape templates. In *2018 International Conference on 3D Vision (3DV)*, pages 672–681. IEEE, 2018. 1, 2
- [14] A. Golovinskiy and T. Funkhouser. Consistent segmentation of 3d models. *Computers & Graphics*, 33(3):262–269, 2009. 2
- [15] T. Groueix, M. Fisher, V. G. Kim, B. C. Russell, and M. Aubry. A papier-mâché approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 216–224, 2018. 3
- [16] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778, 2016. 11
- [17] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015. 6, 11
- [18] R. Hu, M. Savva, and O. van Kaick. Functionality representations and applications for shape analysis. In *Computer Graphics Forum*, volume 37, pages 603–624. Wiley Online Library, 2018. 2
- [19] R. Hu, O. van Kaick, Y. Zheng, and M. Savva. Siggraph asia 2016: course notes directions in shape analysis towards functionality. In *SIGGRAPH ASIA 2016 Courses*, page 8. ACM, 2016. 2
- [20] A. Kaiser, J. A. Ybanez Zepeda, and T. Boubekeur. A survey of simple geometric primitives detection methods for captured 3d data. In *Computer Graphics Forum*. Wiley Online Library, 2018. 2
- [21] E. Kalogerakis, A. Hertzmann, and K. Singh. Learning 3d mesh segmentation and labeling. *ACM Transactions on Graphics (TOG)*, 29(4):102, 2010. 2
- [22] A. Kanazawa, S. Tulsiani, A. A. Efros, and J. Malik. Learning category-specific mesh reconstruction from image collections. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 371–386, 2018. 3
- [23] A. Kanezaki, Y. Matsushita, and Y. Nishida. Rotationnet: Joint object categorization and pose estimation using multi-views from unsupervised viewpoints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5010–5019, 2018. 4
- [24] V. G. Kim, W. Li, N. J. Mitra, S. Chaudhuri, S. DiVerdi, and T. Funkhouser. Learning part-based templates from large collections of 3d shapes. *ACM Transactions on Graphics (TOG)*, 32(4):70, 2013. 2
- [25] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 11
- [26] H. Laga, Y. Guo, H. Tabia, R. B. Fisher, and M. Bennamoun. *3D Shape Analysis: Fundamentals, Theory, and Applications*. John Wiley & Sons, 2018. 2
- [27] V. Léon, V. Itier, N. Bonneel, G. Lavoué, and J.-P. Vandeborre. Semantic correspondence across 3d models for example-based modeling. In *Eurographics Workshop on 3D Object Retrieval 2017 (3DOR 2017)*, 2017. 2
- [28] J. Li, K. Xu, S. Chaudhuri, E. Yumer, H. Zhang, and L. Guibas. Grass: Generative recursive autoencoders for shape structures. *ACM Transactions on Graphics (TOG)*, 36(4):52, 2017. 2
- [29] L. Li, M. Sung, A. Dubrovina, L. Yi, and L. Guibas. Supervised fitting of geometric primitives to 3d point clouds. *arXiv preprint arXiv:1811.08988*, 2018. 2
- [30] Y. Li, H. Su, C. R. Qi, N. Fish, D. Cohen-Or, and L. J. Guibas. Joint embeddings of shapes and images via cnn image purification. *ACM Trans. Graph.*, 34(6):234:1–234:12, Oct. 2015. 6

- [31] Y. Li, X. Wu, Y. Chrysathou, A. Sharf, D. Cohen-Or, and N. J. Mitra. Globfit: Consistently fitting primitives by discovering global relations. *ACM Transactions on Graphics (TOG)*, 30(4):52, 2011. 2
- [32] W. E. Lorensen and H. E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *Proceedings of the 14th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '87*, pages 163–169, New York, NY, USA, 1987. ACM. 2
- [33] D. Maturana and S. Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 922–928. IEEE, 2015. 4
- [34] L. M. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger. Occupancy networks: Learning 3d reconstruction in function space. *CoRR*, abs/1812.03828, 2018. 3, 4
- [35] N. J. Mitra, M. Wand, H. Zhang, D. Cohen-Or, V. Kim, and Q.-X. Huang. Structure-aware shape processing. In *ACM SIGGRAPH 2014 Courses*, page 13. ACM, 2014. 2
- [36] S. Muraki. Volumetric shape description of range data using blobby model. *ACM SIGGRAPH computer graphics*, 25(4):227–235, 1991. 2
- [37] K. Museth. Vdb: High-resolution sparse volumes with dynamic topology. *ACM Trans. Graph.*, 32(3):27:1–27:22, July 2013. 4
- [38] Y. Ohtake, A. Belyaev, M. Alexa, G. Turk, and H.-P. Seidel. *Multi-level partition of unity implicits*, volume 22. ACM, 2003. 1, 3
- [39] M. Ovsjanikov, W. Li, L. Guibas, and N. J. Mitra. Exploration of continuous variability in collections of 3d shapes. *ACM Transactions on Graphics (TOG)*, 30(4):33, 2011. 2
- [40] J. J. Park, P. Florence, J. Straub, R. A. Newcombe, and S. Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. *CoRR*, abs/1901.05103, 2019. 3
- [41] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *arXiv preprint arXiv:1612.00593*, 2016. 4
- [42] A. Ricci. A constructive geometry for computer graphics. *The Computer Journal*, 16(2):157–160, 1973. 2
- [43] L. Roberts. *Machine Perception of Three-Dimensional Solids*. 01 1963. 1, 2
- [44] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. *CoRR*, abs/1505.04597, 2015. 7, 12
- [45] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 11
- [46] G. Sharma, R. Goyal, D. Liu, E. Kalogerakis, and S. Maji. Csgnet: Neural shape parser for constructive solid geometry. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5515–5523, 2018. 2
- [47] V. Sitzmann, J. Thies, F. Heide, M. Nießner, G. Wetzstein, and M. Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. *arXiv*, 2018. 3
- [48] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese. Deep metric learning via lifted structured feature embedding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 8
- [49] H. Su, S. Maji, E. Kalogerakis, and E. G. Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proc. ICCV*, 2015. 4
- [50] M. Tatarchenko, A. Dosovitskiy, and T. Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2088–2096, 2017. 3
- [51] S. Tulsiani, H. Su, L. J. Guibas, A. A. Efros, and J. Malik. Learning shape abstractions by assembling volumetric primitives. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2635–2643, 2017. 2, 3, 6, 8
- [52] L. van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008. 1, 6
- [53] N. Wang, Y. Zhang, Z. Li, Y. Fu, W. Liu, and Y.-G. Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 52–67, 2018. 3, 7, 8
- [54] L. Wei, Q. Huang, D. Ceylan, E. Vouga, and H. Li. Dense human body correspondences using convolutional networks. In *Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 7
- [55] J. Wu, C. Zhang, T. Xue, B. Freeman, and J. Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *Advances in neural information processing systems*, pages 82–90, 2016. 3
- [56] 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. 3
- [57] G. Wyvill, C. McPheeters, and B. Wyvill. Data structure for soft objects. In *Advanced Computer Graphics*, pages 113–128. Springer, 1986. 1, 2
- [58] K. Xu, V. G. Kim, Q. Huang, and E. Kalogerakis. Data-driven shape analysis and processing. In *Computer Graphics Forum*, volume 36, pages 101–132. Wiley Online Library, 2017. 2
- [59] M. E. Yumer and L. B. Kara. Co-abstraction of shape collections. *ACM Transactions on Graphics (TOG)*, 31(6):166, 2012. 2
- [60] Y. Zheng, D. Cohen-Or, M. Averkiou, and N. J. Mitra. Recurring part arrangements in shape collections. In *Computer Graphics Forum*, volume 33, pages 115–124. Wiley Online Library, 2014. 2
- [61] C. Zou, E. Yumer, J. Yang, D. Ceylan, and D. Hoiem. 3d-prnn: Generating shape primitives with recurrent neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 900–909, 2017. 2

A. Additional Algorithmic Details

Hyperparameters: We first detail the hyperparameters associated with training our 3D \rightarrow template network in Table 3. The goal of our architecture design was to enable encoding 3D shape into our representation in a robust and generalizable way. As a result, the architecture is a sequence of convolutional layers followed by a sequence of fully connected layers, and was trained using Adam [25]. We experimented with more complex architectures, such as ResNet [16], but found them to be slow to converge and to generalize poorly (likely due to lack of sufficient training data to train an architecture like ResNet from scratch).

Single View Reconstruction Architecture: We next describe the architecture used for our RGB \rightarrow template experiments (Table 4). In this application, we use ResNet V2 50 to encode images. We find that using this network was more effective than the architecture used for encoding depth images, likely for at least two reasons. First, it could leverage initial weights pretrained on ImageNet classification [45] (all layers except for the final two FC layers were pretrained). Second, it had larger capacity to utilize the larger RGB dataset – while there were approximately 30,000 shapes in the 3DR2-N2 [10] training split of ShapeNet [8], there are over 700,000 images in the 3DR2-N2 RGB render training split. This is because there are 24 RGB renders of each ShapeNet shape in the 3DR2-N2

Name	Value
α	100.0
β	10.0
w_U	1.0
w_S	0.1
w_a	10/3
w_b	0.01
ℓ	-0.07
Batch Size	8
Learning Rate	5×10^{-5}
Adam β_1	0.9
Adam β_2	0.999
Uniform Samples	3000
Near Surface Samples	3000
Conv. Width(s)	3
Conv. Stride(s)	2
Conv. Depths	16, 32, 64, 128, 128
FC Widths	1024, 2048, 2048, 7*N
Nonlinearity	Leaky ReLU
Input Resolution	137x137x20

Table 3. Hyperparameters and optimization details for the early fusion depth image network. This network is composed of 5 convolutional layers followed by 4 fully connected layers.

Name	Value
Batch Size	128
Learning Rate	5×10^{-5}
Backbone Network	ResNet V2 50
Pretraining	ImageNet
Finetuning	End-to-end
FC Widths	2048, 7*N
Image Native Resolution	137x137 (3D-R2N2)
Network Input Resolution	224x224
Adam β_1	0.9
Adam β_2	0.999

Table 4. Hyperparameters and optimization details for the RGB \rightarrow 3D network. This network is composed of ResNet V2 50, up to and including the average pooling layer, followed by two fully connected layers that map to the 3D representation.

set. This acts as a form of data augmentation, which can be leveraged by a larger network.

Note that when training our RGB network, we used network distillation [17] rather than the loss functions described in the main paper. In particular, the only loss for the RGB network was a supervised L^2 regression loss to the “ground truth” template generated by our depth network for each training example.

B. Additional Dataset Details

For all experiments, we used the Shapenet [8] 80%-20% train-test split provided by 3D-R2N2 [10]. Because this split does not include a validation set, we further randomly split train into a 75%-5% train-val split. All networks were trained on train only (not train+val). Val was used to choose hyperparameters, although some experiments in the paper are primarily concerned with train performance.

In addition to the 3D-R2N2 training set, we generated and trained on a larger, more diverse set of 2 million RGB renders of ShapeNet. The renders have a higher native resolution (256x256, compared to 137x137). We also added additional data augmentation during training, varying the brightness, contrast, hue, and saturation randomly. We qualitatively found this to significantly improve domain transfer from synthetic renders to real images. However, for consistency, all results reported for the RGB \rightarrow template network were trained on 3D-R2N2 renders. For the semantic segmentation network, which was never tested on 3D-R2N2, all results were trained on the larger dataset.

C. Runtime Analysis

Inference time: Because our depth \rightarrow template network is a small feed-forward CNN, inference for a template given depth renderings of a mesh is very fast. Using a GTX

Name	Value
Batch Size	8
Learning Rate	1×10^{-4}
Encoder Network	ResNet V2 50
Decoder Network	U-Net based on [44]
Pretraining	None
Input Resolution	256x256
Output Resolution	256x256
Adam β_1	0.9
Adam β_2	0.999

Table 5. Hyperparameters for the 2D semantic segmentation model.

1080 with a batch size of 1, mean network inference time is 1.11ms.

Training time: Since inference is so quick, the performance bottleneck shifts to pre/post processing, such as rendering depth images of the input mesh. The time required to do this is highly dependent on the complexity of the input mesh (although for most of ShapeNet the rasterization process was bottlenecked by the cost of writing the output images to disk).

Surface reconstruction: Though not the goal of our work, structured implicit functions provide a direct way to extract a surface reconstruction (isosurface) from template parameters after inference. Several methods are available with different accuracy vs. time trade-offs.

One possibility is to render the isosurface of the representation directly using ray marching. One benefit of this approach is that the resulting image will accurately reflect the predicted isosurface and its normals. However, achieving interactive framerates with this approach is difficult.

A second alternative is to using marching cubes to extract the isosurface as a mesh. This has the advantage of quick rasterization and it makes comparisons to ground truth meshes easier, so it is the technique we apply in this paper. For our surface extraction implementation, the primary computational bottleneck is to sample $F(\mathbf{x}, \Theta)$ at the marching cubes locations. The simplest implementation would iterate over the elements of the marching cubes grid and compute $F(\mathbf{x}, \Theta)$ directly at each location. However, it is much more efficient to take advantage of the local nature of $f_i(\mathbf{x}, \theta_i)$. In particular, we set a minimum “influence” epsilon for $\frac{f_i(\mathbf{x}, \theta_i)}{c_i}$. This function is the component of $f_i(\mathbf{x}, \theta_i)$ that falls off with distance and is always in $[0, 1]$. In practice we set $\epsilon = 10^{-3}$. Then, we iterate over each shape element, and add its marginal contribution to the marching cubes volume by iterating only over the voxels where it has nontrivial contribution. With this algorithm, mesh extraction takes approximately 5.96 seconds at 256^3 on the CPU with one thread. GPU acceleration for this

task is also possible, which could result in further runtime decreases.

We also postprocess the extracted isosurface by removing connected components below a trivial size threshold. This is helpful because one common artifact is for a function $f_i(\mathbf{x}, \theta_i)$ to be almost off, but still cross ℓ within a small region of space. As the marching cubes resolution increases, smaller and small volumes can be revealed, so this post processing step is useful. A surface area threshold of 0.005 for connected components was used for the interpolation video.