# Parsing Geometry Using Structure-Aware Shape Templates

Vignesh Ganapathi-Subramanian
Stanford University
vigansub@stanford.edu

Olga Diamanti
Autodesk Inc., Stanford University
diamanti@stanford.edu

Soeren Pirk
Stanford University
pirk@google.com

Chengcheng Tang
Stanford University
tangcc@stanford.edu

Matthias Nießner
TU Munich, Stanford University
niessner@tum.de

Leonidas J. Guibas
Stanford University
guibas@cs.stanford.edu

## Abstract

*Real-life man-made objects often exhibit strong and easily-identifiable structure, as a direct result of their design or their intended functionality. Structure typically appears in the form of individual parts and their arrangement. Knowing about object structure can be an important cue for object recognition and scene understanding - a key goal for various AR and robotics applications. However, commodity RGB-D sensors used in these scenarios only produce raw, unorganized point clouds, without structural information about the captured scene. Moreover, the generated data is commonly partial and susceptible to artifacts and noise, which makes inferring the structure of scanned objects challenging. In this paper, we organize large shape collections into parameterized shape templates to capture the underlying structure of the objects. The templates allow us to transfer the structural information onto new objects and incomplete scans. We employ a deep neural network that matches the partial scan with one of the shape templates, then match and fit it to complete and detailed models from the collection. This allows us to faithfully label its parts and to guide the reconstruction of the scanned object. We showcase the effectiveness of our method by comparing it to other state-of-the-art approaches.*

## 1. Introduction

In all their variability, real-life man-made objects are often designed with certain basic principles in mind, relating to their target functionality, affordances, physical and material constraints, or even aesthetics. As a result, most objects can be described by common structural forms or patterns. For example, chairs often follow a model of "four-legs", with/without armrests, "S-"shaped, or swivel. Object structure can be modeled by identifying the most common part layouts, spatial interrelations and part symmetries.
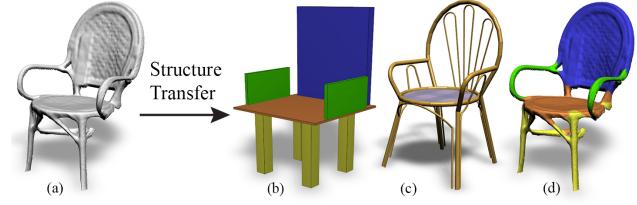


Figure 1: Structure transfer with templates. Given a partial scan (a), we detect and fit a matching structural template (b). A similar complete shape from a collection, recovered and deformed to fit the scan (c). (d) Semantic part labels transferred onto the scan.

Associating an unknown object with a particular structural pattern can help with understanding the object's properties and purpose. In particular, structure can help in cases where part of the object geometry is missing, due to the object being only partially observed. This is a common scenario in 3D sensing pipelines, e.g. those arising in augmented reality (AR) and robotics, or used for data acquisition for graphics/visual effects applications, where a scene is reconstructed for further editing. 3D sensing pipelines typically rely on scanning a scene through RGB-D sensors, which commonly produce incomplete, noisy, and unorganized point scans. The resulting reconstructed geometry often is difficult to parse and may fail to meet the requirements of applications that use more precise representations of shapes, e.g. CAD models. Faithfully inferring object structure and using it to recover a corresponding complete surface mesh from partial sensor data remains one of the fundamental problems in shape reconstruction research.

We introduce a pipeline for transferring structural information from large shape collections onto unknown scanned shapes. We model object structure via a small number of hand-crafted templates, to capture the structural patterns exhibited by shapes of a particular category. Our hand-crafted templates are designed to represent an abstraction of the shape. This is similar to how designers model shapes, often starting with a rough structure in mind and based on the

desired utility – only afterwards are the details carved out. Manually providing the templates is meant to emulate this process. Our templates are abstract, lightweight, yet meaningful enough to account for thousands of shapes in various forms. Each template consists of axis-aligned boxes in a specific spatial arrangement, providing knowledge about the structural organization of the shape and the relationship of its parts. We parameterize the templates so as to fit them, via state-of-the-art optimization, to a particular shape.

We then leverage the structural information encoded by the templates for the entire shape collection to learn their structure. This allows us to identify partial shapes obtained through scanning. To address this problem of inferring the structure of a shape from its partial geometry, we employ a deep neural network trained on partial views of shapes. The network is able to detect shape structure. Using this, one can identify the scanned object by retrieving the closest shape from the collection and fitting it to the scan. Additional applications of the templates include part labels for partial scans, and recovery of a fully meshed CAD model to replace the scan in scene modeling scenarios (Fig. 1).

## 2. Related Work

Due to the complexity and variability of objects, efforts on representing and understanding shapes focus on their structure [21], their features and similarities [1], the semantic meaning of individual parts [35], and even the creative process of shape modeling [5].

**Structure-aware Representations.** It has been recognized that the structural organization of shapes plays an eminent role in modeling and reconstructing shapes. Existing approaches focus on identifying shape parameters [10], relations [18, 16], symmetries [24, 36], correspondences [22], combinatorial variations [9, 2], or co-variations [34]. More recently, Schulz et al. [28] show that discrete and continuous shape features can be represented as low-dimensional manifolds and then used to retrieve individual shapes from large shape collections. All these techniques identify and represent shape structure, however, they do not aim at reconstructing or replacing partial scans.

**Shape Templates and Part-based Models.** Shape templates have proven to be an effective tool for inferring higher-level knowledge of shapes, not only in the context of 3D geometry [23], but also for various image processing tasks [8, 7, 20]. Kim et al. [19] employ part-based templates to encode deformations and fit object labels to scanned point clouds. Kalogerakis et al. [16] learn distributions of parts to encode part placements, Kim et al. [17] propose a fully automatic approach for inferring probabilistic part-based templates to capture style and variation of shapes in large model collections. Unlike the existing approaches, we use shape templates to organize large shape repositories and to learn the structure of shapes and their parts.
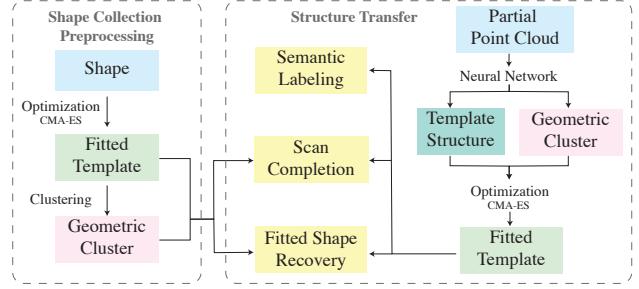


Figure 2: Overview: we propose a two-stage framework for parsing shapes with the end goal of retrieving structure from partial point clouds. In the first stage (shape collection preprocessing - left), we organize the shapes in our collection according to their structure, using templates. In the second stage (structure transfer to partial shapes - right), we train a deep neural network to retrieve a proper geometry-based cluster and a suitable structure template. We further refine the template parameters through optimization.

**Shape Reconstruction.** An especially difficult problem is the reconstruction of shapes from unstructured or incomplete data, such as partial meshes or point clouds. Symmetry-driven reconstruction leverages symmetric properties in the scanned data to complete occluded regions [24, 3]. Sipiran et al. [31] rely on local features and fit surface functions to incomplete point clouds. Shen et al. [30] propose a bottom up structure recovery pipeline that allows to replace incomplete point scans of man-made objects by aligning them to existing shapes in a repository. Sung et al. [32] use 3D shape collections and exploit symmetry relations to infer global structure and complete point clouds with substantial occlusion. However, they concentrate on learning the distributions of shape part features, and use the collection to retrieve only individual shape parts, not entire shapes to match the partial point cloud.

## 3. Overview

Our framework works in two stages (Fig. 2). In the first step, we fit a set of pre-defined shape templates to shapes of a shape collection. This is necessary, as shape repositories commonly do not provide information about the structure of shapes, but instead only store surface meshes of the models. Fitting templates allows us to infer the structural organization of the shape collection. Moreover, we cluster shapes according to their template parameterization. This provides a geometric organization of the shapes, which can be leveraged to perform approximate shape retrieval from the collection. In the second step, we employ the template organization imposed on the collection above to infer shape structure for unseen shapes with partial geometry. We use the inferred structure to retrieve and fit known shapes from the database to the partial scan, which can directly be used as a proxy shape for scanned geometry. Additionally, the structure allows to identify and annotate shape parts and to

reconstruct the object in a structure-aware manner.

# 4. Shape Collection Preprocessing

In this section we describe the first stage of our pipeline, which involves organizing the shapes in the collection $\mathcal{S}$ using our primitive-based shape templates. They provide a structure-aware shape representation. We assume that all the models in $\mathcal{S}$ are pre-aligned, which makes defining a common frame of reference possible.

## 4.1. Structure-Aware Templates

We first define a set of structural templates. Each template captures a particular structural "mode" frequently found among shapes in that family – in the "chair" example, one template might capture the common four-legged pattern, while another the swivel structure. We fit these templates to all shapes of that family in the collection, and choose the most appropriate template for each shape.

**Template Definition.** A *template* is a collection of deformable axis-aligned box primitives and connectors, forming a tree graph structure. Each box represents a structural part of a shape. Specifically, a template consists of:

**Nodes.** Graph nodes (vertices) are associated with the box primitives in the template. Each box is described by 6 *parameters*, namely its position and size in 3D.

**Connectors.** Template graph edges "glue" together adjacent boxes, structurally attaching them and constraining their dimensional parameters. Connectors also define the relative positions of child boxes to parent boxes. We only model attachment at right angles.

**Groups.** Symmetric relationships are often present among structural parts of shapes, e.g. between the four legs of a four-legged chair. In terms of our templates, symmetries are modeled by symmetric structural constraints between graph nodes, which require that the parts (in this example, the four leg boxes) are identical. This symmetry meta-information is encoded via groups in the graph.

We pre-design a set of $N_T = 21$ templates in total. Each shape family is associated with a subset of these templates, but a single template can be shared by more than one shape family. For example, tables and desks often have similar structure, which is reflected by them sharing one of their templates. The graph structure of the templates stays fixed throughout our pipeline, but we can tune a template to a particular shape by finding an optimal configuration of template box parameters. The fitted template is then a structure-aware representation of the shape, with semantic information about the shape parts encoded as meta-information.

## 4.2. Template Selection and Fitting

We are now given a shape $S$ of a particular family, in the form of an unstructured point cloud sampled from a
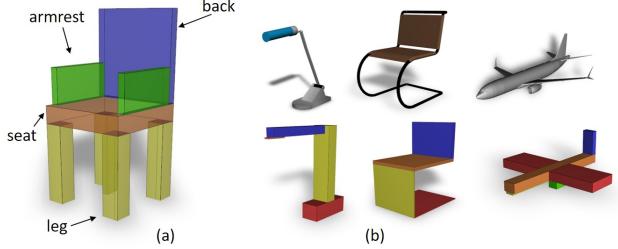


Figure 3: (a) A shape template with semantic information embedded in its parts (b) Top - various shapes from the collection: a lamp, a chair with folding legs, and an airplane. Bottom - The fitted templates, with parameters that best approximate those shapes.

database shape. The aim is to find which template structure $T = \mathcal{T}(S)$ best approximates the structure of $S$, and compute values for its parameters $\mathcal{P}(S, T)$ (box sizes and locations) to fit the shape's geometry. We proceed in two stages: first, we fit all templates pertaining to the shape's family to the shape, e.g. chair templates to a chair shape. Then, we select the best-fitting template from that set.

### 4.2.1 Template Fitting

We aim to fit a template structure $T$ to an input point cloud $S$. The template consists of boxes $B_i$, $i = 1, 2, ..., N$. For simplicity, in the following discussion we denote both the template and its shape-dependent parameters $\mathcal{P}(S, T)$ as $T$. The optimal parameter values are found by solving an optimization problem $T^* = \underset{T}{argmin}\, E_{\text{total}}(T) = \sum_i \lambda_i E_i(T)$ where $i$ ranges over individual energy terms. The various terms in $E_{\text{total}}$ encourage a close match between the "box-like" template shape and the input point cloud at the optimum. We used $\lambda_{\text{proj}} = 0.3$, $\lambda_{\text{bbox}} = 1$, $\lambda_{\text{min}} = 0.8$, $\lambda_{\text{dis-ent}} = 0.4$ in all our experiments. The individual energy terms are detailed below. A qualitative evaluation of various energy terms is discussed in the supplementary material.

**Projection:** The sum of distances from all points in the point cloud $S$ to the template geometry: $E_{\text{proj}}(T) = \sum_{\mathbf{p} \in S} \min_{j=1,2,...,N} d(\mathbf{p}, B_j)$, where $d(\mathbf{p}, B_j)$ is the minimum projection distance from point $\mathbf{p} \in \mathbb{R}^3$ to box $B_j$ in the template. The projection term ensures that the optimization produces a well-fitting template.

**Bounding Box:** The difference in size (3D volumes) between the bounding box of the point cloud $S$ and that of the template $T$: $E_{\text{bbox}}(T) = |Vol(T_{\text{bbox}}) - Vol(S_{\text{bbox}})|$

**Minimalism:** The total volume of the template in space: $E_{\text{min}}(T) = \sum_{B_i \in T} Vol(B_i)$ With this term, the thinnest possible set of boxes is fitted to the point cloud, ensuring that the template geometry is no bigger than need be. While this term overlaps with the bounding-box energy, we found that it promotes faster optimization convergence.

**Disentanglement:** The amount of overlap between boxes: $E_{\text{dis-ent}}(T) = \sum_{B_i \in T} \sum_{B_j \in T, B_j \neq B_i} Vol(B_i \cap B_j)$ This term requires that template boxes don't spatially obstruct each other, since they are meant to capture distinct semantic shape parts.

### 4.2.2 Optimization and Template Selection

The energy $E_{\text{total}}(T)$ is highly non-convex, requiring an optimization scheme that avoids local minima. Evolutionary optimization strategies are particularly appropriate for this task [13], since they continuously adapt the step size based on previous observations and utilize a randomization factor to probe large energy landscapes. In addition, they are more easily parallelizable compared to classical gradient-based techniques. We employ the Covariance Matrix Analysis - Evolutionary Strategy (CMA-ES) [12], which uses correlations between different variables to determine the step size.

Given the non-convexity and in order to aid convergence, we initialize the optimization based on solutions of previously successful optimizations for other shapes of the same family. We eventually choose the optimal set of parameter values across all runs as the best fit for a given template. Once all candidate templates have been fitted to a given shape, we select the best-fitting one $T^*$ as the one that minimizes $E_{\text{total}}(T)$ over all candidate templates. A caveat to performing template selection when two similar templates are used and also various convergence timings are discussed in the supplementary material.

### 4.3. Structural Shape Clustering

We further use the templates to cluster the shapes in the collection into groups of similar structure and rough dimensions. We first group the shapes in the collection according to their best fitting template structure. Then, we use the vectors of template parameters (box dimensions) to further divide the shapes of each group into clusters of shapes of similar part dimensions via k-means. We use 10 clusters per template ($N_C = 210$ clusters total). Each cluster is then associated with a specific template, but can contain shapes from different families, since the same template can be shared by more than one family (e.g. tables and desks). The clusters will be used to inform the structure identification of partially scanned shapes.

## 5. Structure Transfer to Partial Shapes

Having used the templates to organize the data in our shape collection in terms of their structure, we now transfer structure to partial scans of objects not present in the shape collection. We use the structural information to better understand the partial object's shape, and structure it by assigning and fitting it to one of the known structural patterns found in the collection.
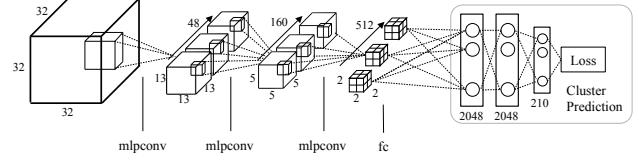


Figure 4: Classification Network, inspired by [27]. Multiple mlpconv layers, which are powerful local feature extractors, are followed by a set of fully connected layers. This network learns and shape cluster from the input signed distance field of a partial point scan, and from that, the best structural template.

### 5.1. Inferring Structure from Partial Scans

The input to the structure inference stage is a partial point scan $S_{\text{partial}}$. The output is one of the templates, fitted to $S_{\text{partial}}$; this process imposes structure on the partial point cloud. We train a deep neural network that (indirectly) assigns the partial shape to a particular shape template; it predicts structure from partial geometry. Trained with simulated partial views, the network learns to ignore various artifacts and noise commonly found in real RGB-D scans. Once the network selects the template pattern, we further optimize template parameters to fit the scan ( Section 4.2.1).

Note that the template assignment, on its own, is not sufficient to identify a shape, since a particular template graph can be shared by multiple shape families. Not knowing the shape family has the additional negative side-effect that we cannot intelligently initialize the fitting CMA-ES optimization as described in Section 4.2.2, which can adversely affect performance both in terms of runtime and quality of results. Thus, instead of having the network directly predict the template structure, we train it to predict the *shape cluster* instead; since each cluster is associated with a particular template structure, the network also indirectly predicts the template. The predicted cluster also provides, via its shapes, a set of candidate template fits for that particular template pattern, which can be used as initialization to get the optimal fit of the template pattern to the partial scan. Thus, the network predicts, from raw partial geometry, both a rough structure pattern and part dimensions for the partial scan.

### 5.1.1 Classification Network

The input to the network is the signed distance field $d_{\text{partial}}$ of the partial point scan $S_{\text{partial}}$. The output is a vector of length $N_C$, encoding the probabilities with which each shape cluster corresponds to the scan. The network architecture (Fig. 4) uses mlpconv layers[27], shown to be powerful local patch feature extractors. They are followed by a number of fully connected layers and dropout layers of factor 0.2. We use cross-entropy [11] as the loss function and perform data augmentation [27] to reduce overfitting.

Figure 5: Structure-aware shape manipulation. Left to right: source shape and fitted template, target shape and fitted template, deformed version of the source shape to match the target. The rightmost heat map indicates the normalized per-vertex magnitude of deformation induced on the source shape.

### 5.1.2 Partial Shape Identification

The most likely shape clusters, as predicted by the network, indicate the most likely template structures for the partial point cloud $S_{\text{partial}}$. We fit templates corresponding to the $k = 3$ most likely clusters to the partial shape (Section 4.2.1) and pick the best-fitting template $T_{\text{partial}}$ based on fitting error. This fully structures the partial shape and identifies its family. We initialize the optimization for each template by averaging the best-fitting parameter values for the shapes in the cluster where the template came from. The optimization also produces the optimal template parameters $\mathcal{P}(S_{\text{partial}}, T_{\text{partial}})$ aligning the template to partial geometry.

### 5.2. Structure-Aware Partial Shape Recovery

The identification of a partial shape, via its fitted template and family, enables *shape recovery*: we can retrieve, among all shapes in the collection, one that fits the partial point cloud. This is useful for scene understanding scenarios, in AR or robotics, where a partial point cloud can be mapped to (and possibly replaced by) an already known shape, or for editing a scanned scene for CG applications. Since we cannot directly use shape geometry to detect the most similar shape as our input is partial, we use fitted templates to provide a rough proxy for geometric and structural similarity, via their box dimensions and locations. We thus look for a shape in the collection that matches $S_{\text{partial}}$ in terms of its parametric template fit. We search in the $k = 3$ most likely clusters as provided by the network. The output of this stage is a shape from the collection $S_{\text{source}}$, along with the optimal parameter values $\mathcal{P}(S_{\text{source}}, T_{\text{partial}})$ that fit $S_{\text{source}}$ to the template of the partial shape, $T_{\text{partial}}$.

### 5.2.1 Fitting the known shape to the scan

Our templates provide an intuitive way of simultaneously manipulating all shapes sharing the template structure: in this case, we can deform the collection shape $S_{\text{source}}$ to
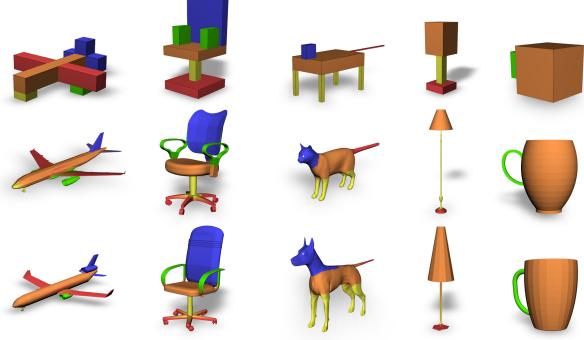


Figure 6: Semantic labeling for airplanes, chairs, animals, lamps and mugs, as produced by the box-template parameters. The parts are color-coded based on their tags: *airplanes*-{(orange, body), (red, wings), (blue, tail), (green, engine), (yellow, front wheel)}, *chairs*-{(orange, seat), (red, base), (yellow, swivel), (green, armrests), (blue, backrest)}, *animals*-{(blue, head), (orange, torso), (red, tail), (yellow, legs)}, *lamps*-{(orange, bulb), (yellow, stem), (red, base)}, *mugs*-{(orange, cup),(green, handle)}.



Figure 7: Semantic labeling for partial point scans of shapes, using our pipeline for structure transfer.

match the partial scan $S_{\text{partial}}$. This will recover a complete shape matching the scan. Since finding $S_{\text{source}}$ takes into account the template fit, the amount of distortion that this deformation process induces to $S_{\text{source}}$ is typically minimal.

**Template-based Deformation**  Assume we are given two shapes $S_1$, $S_2$ with compatible structure; namely, their best-fitting templates (Section 4.2) have the same graph structure, but different parameter values. We denote the two parameterized fitted templates for the two models by $B_1$ and $B_2$, both with $n$ boxes. The goal now is to "morph" $S_1$ to $S_2$. Our simple approach is inspired by traditional skinning techniques [29, 14, 15]: since there is a box-to-box correspondence between $B_1$ and $B_2$, we define the transformation of $S_1$ to $S_2$ via a weighted sum of the affine transformations $\{\mathcal{A}_i(.)\}$ that map each individual box $B_{1,i}$ to its corresponding box $B_{2,i}$, for $i = 1, 2, ..., n$. The parameters of $B_{1,i}$ are its center $\mathbf{c}_1 = [c_1^x, c_1^y, c_1^z]^T$ and its dimensions $\mathbf{l}_1 = [l_1^x, l_1^y, l_1^z]^T$ – similarly for $B_{2,i}$. Then, the affine transformation mapping each point $\mathbf{p}_1$ on $B_{1,i}$ to its corresponding point $\mathbf{p}_2$ on $B_{2,i}$ is given by $\mathbf{p}_2 = \mathbf{c}_2 + \mathbf{R}_{12}(\mathbf{p}_1 - \mathbf{c}_1)$ where $\mathbf{R}_{12} = \text{diag}\left(\frac{l_2^x}{l_1^x}, \frac{l_2^y}{l_1^y}, \frac{l_2^z}{l_1^z}\right)$. Using the individual box

| Category | #Shapes | PointNet | SyncSpecCNN | Ours |
|----------|---------|----------|-------------|------|
| Chair | 3746 | 89.6 | **90.24** | 87.37 |
| Table | 4520 | 80.6 | 82.13 | **88.62** |
| Cup | 184 | 93.0 | 92.73 | **94.06** |
| Lamp | 1547 | 80.8 | **84.65** | 78.03 |
| Airplane | 2690 | **83.4** | 81.55 | 76.71 |

Table 1: Intersection-over-Union (IoU) percentages for our part labeling, evaluated against the ground truth and compared to [26] and [38]. The highest IoU value for each category is in bold.

| Task | Accuracy | | |
|------|----------|--------|--------|
| | Best | Best 2 | Best 3 |
| Cluster Classification - Training | 89% | 94% | 98% |
| Cluster Classification - Testing | 73 % | 83 % | 92 % |

Table 2: Performance of the dual-classification network for template and cluster classification on partial virtual scans of ShapeNet objects. The three columns report the percentage of cases where the correct fit (cluster or template) was found among the 1, 2 or 3 most likely fits according to the net's output.

transformations, any given point $\mathbf{q}_1$ on $S_1$ is mapped to a point $\mathbf{q}_2$ on $S_2$ via: $\mathbf{q}_2 = \mathcal{A}(\mathbf{q}_1) = \frac{\sum_i w_i(\mathbf{q}_1)\mathcal{A}_i(\mathbf{q}_1)}{\sum_i w_i(\mathbf{q}_1)}$. This yields a continuous full shape transformation. We choose the weights to be in inverse relation to the distance between the point and the boxes, to ensure that any one box only locally affects the points inside and around it. Additionally, we want the weights to be smooth, to preserve smoothness of the underlying geometry. We use $w_i(\mathbf{q}) = \exp\left(-d(\mathbf{q}, i)^2\right)$, where $d(\mathbf{q}, i)$ denotes the distance between point $\mathbf{q}$ and the $i$-th box in the template. After mapping all points on $S_1$ using the process above, we perform a global scaling so that the final deformed shape $\tilde{S}_1$ lies in the same global bounding box as $B_2$ and thus matches its proportions, but preserves the shape details from $S_1$. Note that $S_1$ and $S_2$ can be represented by either meshes or point clouds at this stage. In the former case, we transform the vertices of $S_1$ and keep the connectivity the same.

**Fitting source to partial** In order to recover the best fitting shape for our partial scan, we simply apply the process above with $S_1 = S_{\text{source}}$ and $S_2 = S_{\text{partial}}$. The output $\tilde{S}_{\text{source}}$ is the recovered model.

# 6. Results and Discussion

The pipeline has been implemented in C++ – we plan to make the code freely available online. Our shape collection is a subset of the ShapeNet repository [4]. We used shapes from 10 different categories: monitors, cups, tables, chairs, benches, desks, dressers, sofas, airplanes, and lamps.
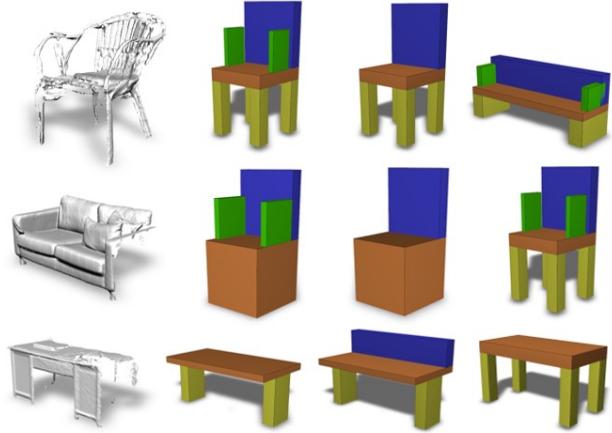


Figure 8: The three most likely templates, as determined by the neural network for given partial scans. Left to right: input partial scan and possible structure templates (pre-parameter fitting), in decreasing order of likelihood.

## 6.1. Template Fitting and Selection for Database Shapes

We fit templates running the CMA-ES [12] algorithm for four different initializations. Some examples of fitting template parameters to shapes in the collection $\mathcal{S}$ (Section 4) are shown in Fig. 3. The success of this process is critical for various other results, eg. shape manipulation (Fig. 5) and part labeling (Fig. 6), and thus fitting results are implicitly showcased as an intermediate step in these figures.

## 6.2. Template-Based Deformation

To evaluate our template-based deformation (Section 5.2.1), we show some results with complete shapes in Fig. 5. We show the source $S_1$ and target $\tilde{S}_1$ shapes, their fitted templates, and the deformed version of $S_1$ that aligns to $\tilde{S}_1$. Since all shapes are scaled to fit the figure, we measure the magnitude of induced deformation by via the normalized Euclidean distance between each vertex in $S_1$ and its mapped location in $\tilde{S}_1$, and show it as a heat map. The box-based transformations help individually adapt, and globally align, the individual semantic parts, without drastically losing their individuality (while the overall sizes of seats in Fig. 5 match, they do not deform to align with one another). This enables generating variations of shapes in a structure-preserving way, which could be interesting for editing or morphing applications; this is however not a focus of this paper.

## 6.3. Deep Network Training and Output

Table 2 shows the accuracy of cluster classification achieved by the classification network for partial scans of shapes. The network operates on 22051 signed distance fields obtained from virtually generated partial scans of

Figure 9: Partial shape recovery on the real-data reconstruction dataset [27], for 4 different categories, bench, monitor, mug, desk. Top row: partial scans of entirely new objects (i.e. not present in the shape collection). Middle: optimization-based template fitting to the partial scan, after classification to template structure and geometric cluster. Bottom: recovered shapes from the collection, deformed to fit the partial scan.

shapes from the ShapeNet repository, across the 8 different shape categories. We used 80% of this data for training and the remaining for testing. As is also evident in Table 2, using the net to select the *single* best cluster for a partially scanned shape may be a suboptimal strategy, since the classification accuracy is not always high. Our strategy of exploring the best three cluster predictions ensures that we are operating at very high classification accuracy. This, in return, ensures that the template fitting is initialized with a much better estimate. The choice of a number of top clusters thus provides a trade-off between quality of initialization and post-processing time needed to perform the optimization rounds.

Qualitative results of the network output on new partial shape scans are shown in Fig. 8. The figure shows the most likely templates, as established by the clusters, which at this stage encode the structure and also provide a rough initialization for the template fitting. (Section 5.2).

## 6.4. Partial Shape Recovery

Recovering shapes from partial point scans, using the network output and the identification/fitting process of Sections 5.1.2 and 5.2 is shown in Fig. 9. Our template fitting optimization is robust to partiality as well as noise in the point scans. This is highlighted in more detail in the supplementary material. As can be observed in the last two columns of Fig. 9, partial shape recovery on shapes that do not fit the classical box-template structure can be done effectively. In case of the folding chair, the dimensions of the template fit, with a thick backrest and thick legs, place the partial scan into a cluster containing other folding chairs that have the same dimensional properties on the back and

legs. Similarly, though the fit of the S-shaped partial scan does not contain the legs, the dimensions of the remaining parts aid in recovering a similarly shaped object from the database to complete the partial scan. In addition to reconstructing a point cloud by recovering a complete mesh from the collection, we can also utilize the inferred structure to augment the point cloud itself. Since this is not the main focus of the paper, we provide some point cloud completion techniques in the supplementary material with some qualitative results.

**Scene completion.** In Fig. 10, we provide an example of recovering shapes for real-life RGBD scanned indoor scenes. We preprocess the scenes using [25], which detects 3D objects in the scene and annotates the RGBD point cloud with the objects' 3D bounding boxes. We treat the points in each of these bounding boxes as a partial scan, recover a fitted CAD mesh using our structure-aware recovery pipeline, and replace scanned points with the recovered mesh. The retrieved shapes are fairly close to the input shapes despite heavily intersecting bounding boxes. Replacing scanned points in the bounding boxes with retrieved shapes makes the scene less cluttered and would allow for further scene editing. Note that the failure to recover lamp and sofa meshes in the first scene is due to their bounding boxes not being detected by the aforementioned method.

## 6.5. Comparisons.

In Fig. 11(a), we qualitatively compare against the work of Sung et al. [32], which also retrieves box-based parts for an input partial point cloud. In their case, parts are retrieved individually via optimization – in comparison, our

Figure 10: Scene completion with full mesh replacements for identified objects in a partially generated scene. Left to right: RGB images with bounding box annotations on objects, point cloud generated using [25], scene completion provided by our fitting method on the bounding boxes.

method tends to provide more realistic part layouts due to the structure enforced by the templates. Fig. 11(b) provides a qualitative comparison to the end-to-end technique of Dai et al. [6]. Since no structure is available, this technique provides somewhat "blobby" completions ; while these results could be used for object classification tasks or as training datasets, the recovered shape is too rough to be used as as a prototype e.g. in AR/CG scene editing applications. In contrast, our pipeline provides a full CAD mesh, fitted to match the input scan, which can be directly be used in place of the partially scanned shape e.g. for scene modeling purposes.

In the inset figure, we show a comparison to the fit produced by Tulsiani et al. [33]. Even though their result is regressed by a network trained towards box-based fitting of chairs specifically, the



lack of a clearly defined structure among the boxes makes it difficult to correctly fit to the shape. In comparison, our constrained templates recover the missing parts of the shape more accurately. For fairness, we note that Tulsiani et al. [33] do not aim to complete partial shapes.

### 6.6. Semantic Part Labeling

The box-primitives in our shape templates are, by design, associated with real shape parts with a "semantic" meaning – e.g. legs or armrests in a chair template. Fitting a template transfers this information onto a shape and annotates it with meaningful parts. We annotate each point on a source shape $S$ by assigning the index of the box closest to it (by projection distance) in its best-fitting template $T^*(S)$. This can be done both for partial scans (Fig. 7) and for complete shapes, e.g. from the shape collection (Fig. 6).

Table 1 discusses accuracy of part labeling on 5 categories. We consider the labeling by [37] as ground-truth, where shape part labels are obtained in a semi-supervised



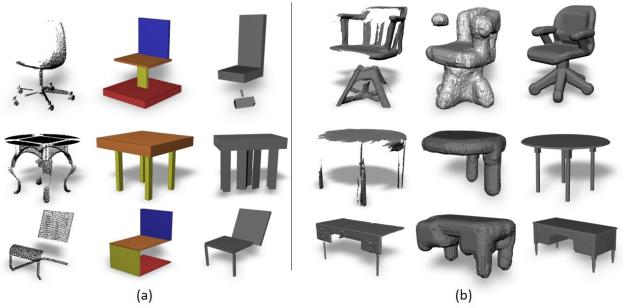(a)                                    (b)

Figure 11: (a) Comparison to Sung et al.[32] Columns, left to right: partial point clouds, our templates fit to the input, box fits generated by Sung et al. to the input. (b) Comparison to Dai et al. Columns, left to right: input partial scans, shape completion obtained by applying the end-to-end technique proposed in Dai et al., our shape retrieval and fitting method applied on the input.

way, using Mechanical Turk verification. We evaluate our technique against the performance of [26] and [38] on the same task. While both these methods use supervision on 70% of the category size in ShapeNet, with an additional 10% of the shapes used in the validation set, we use no external supervision to perform the part labeling.

## 7. Conclusion and Future Work

Obtaining structural information about an object, scanned by commodity hardware into an unstructured partial point cloud, can be key to identifying the object and reasoning about its functionality. We represent structure by a set of pre-designed structural templates, based on simple box-like primitives. We leverage the obtained structural information using a neural network, and show applications of recovering the shape of a partial scan, annotating its structural parts, and applying this to perform scene completion. We provide a single lightweight pipeline that achieves good performance in all these tasks. Our method is unique in that it recovers a full mesh to account for a partial scan of an object. This highlights the value of simple hand-crafted templates, which can abstract away significant geometric detail. That said, automatically inferring the shape templates, and even incorporating different primitives, is the ideal scenario; it remains a difficult unsolved problem, especially when both template parts and inter-part symmetries need to be inferred providing an interesting avenue for future work.

### Acknowledgements

# References

[1] S. Biasotti, A. Cerri, A. Bronstein, and M. Bronstein. Recent trends, applications, and perspectives in 3d shape similarity assessment. *CGF, Online Preprint*, 2015. 2

[2] M. Bokeloh, M. Wand, V. Koltun, and H.-P. Seidel. Pattern-aware shape deformation using sliding dockers. *ACM Trans. Graph.*, 30(6):123:1–123:10, Dec. 2011. 2

[3] M. Bokeloh, M. Wand, and H.-P. Seidel. A connection between partial symmetry and inverse procedural modeling. *ACM Trans. Graph.*, 29(4):104:1–104:10, July 2010. 2

[4] A. X. Chang, T. A. Funkhouser, L. J. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu. Shapenet: An information-rich 3d model repository. *CoRR*, abs/1512.03012, 2015. 6

[5] D. Cohen-Or and H. Zhang. From inspired modeling to creative modeling. *Vis. Comput.*, 32(1):7–14, 2016. 2

[6] A. Dai, C. R. Qi, and M. Nießner. Shape completion using 3d-encoder-predictor cnns and shape synthesis. *http://arxiv.org/abs/1612.00101*, 2016. 8

[7] S. Eslami and C. Williams. A generative model for parts-based object segmentation. In P. Bartlett, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 100–107. 2012. 2

[8] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE Trans. Pattern Anal. Mach. Intell.*, 32(9):1627–1645, Sept. 2010. 2

[9] T. Funkhouser, M. Kazhdan, P. Shilane, P. Min, W. Kiefer, A. Tal, S. Rusinkiewicz, and D. Dobkin. Modeling by example. *ACM Trans. Graph.*, 23(3):652–663, Aug. 2004. 2

[10] R. Gal, O. Sorkine, N. J. Mitra, and D. Cohen-Or. iwires: An analyze-and-edit approach to shape manipulation. *ACM Trans. Graph.*, 28(3):33:1–33:10, July 2009. 2

[11] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org. 4

[12] N. Hansen. The cma evolution strategy: A comparing review. *Towards a new Evolutionary Computation*, 2006. 4, 6

[13] N. Hansen and A. Ostermeier. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation. Morgan Kaufmann, 1996. 4

[14] P. Joshi, M. Meyer, T. DeRose, B. Green, and T. Sanocki. Harmonic coordinates for character articulation. In *ACM SIGGRAPH 2007 Papers*. 5

[15] T. Ju, S. Schaefer, and J. Warren. Mean value coordinates for closed triangular meshes. *ACM TOG*, 24(3), 2005. 5

[16] E. Kalogerakis, S. Chaudhuri, D. Koller, and V. Koltun. A probabilistic model for component-based shape synthesis. *ACM Trans. Graph.*, 31(4):55:1–55:11, July 2012. 2

[17] 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 TOG*, 32(4), 2013. 2

[18] V. G. Kim, W. Li, N. J. Mitra, S. DiVerdi, and T. Funkhouser. Exploring collections of 3d models using fuzzy correspondences. *ACM TOG*, 31(4):54:1–54:11, July 2012. 2

[19] Y. M. Kim, N. J. Mitra, D.-M. Yan, and L. Guibas. Acquiring 3d indoor environments with variability and repetition. *ACM Trans. Graph.*, 31(6):138:1–138:11, Nov. 2012. 2

[20] R. J. López-Sastre, T. Tuytelaars, and S. Savarese. Deformable part models revisited: A performance evaluation for object category pose estimation. In *2011 IEEE ICCV Workshops*, pages 1052–1059, Nov 2011. 2

[21] N. J. Mitra, M. Wand, H. Zhang, D. Cohen-Or, V. G. Kim, and Q.-X. Huang. Structure-Aware Shape Processing. *SIGGRAPH Course notes*, 2014. 2

[22] M. Ovsjanikov, M. Ben-Chen, J. Solomon, A. Butscher, and L. Guibas. Functional maps: A flexible representation of maps between shapes. *ACM Trans. Graph.*, 31(4):30:1–30:11, July 2012. 2

[23] M. Ovsjanikov, W. Li, L. Guibas, and N. J. Mitra. Exploration of continuous variability in collections of 3d shapes. *ACM Trans. Graph.*, 30(4):33:1–33:10, July 2011. 2

[24] M. Pauly, N. J. Mitra, J. Wallner, H. Pottmann, and L. J. Guibas. Discovering structural regularity in 3d geometry. *ACM Trans. Graph.*, 27(3):43:1–43:11, Aug. 2008. 2

[25] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas. Frustum pointnets for 3d object detection from RGB-D data. *CoRR*, abs/1711.08488, 2017. 7, 8

[26] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *CoRR*, abs/1612.00593, 2016. 6, 8

[27] C. R. Qi, H. Su, M. Nießner, A. Dai, M. Yan, and L. J. Guibas. Volumetric and multi-view cnns for object classification on 3d data. *CoRR*, abs/1604.03265, 2016. 4, 7

[28] A. Schulz, A. Shamir, I. Baran, D. I. W. Levin, P. Sitthi-Amorn, and W. Matusik. Retrieval on parametric shape collections. *ACM Transactions on Graphics*, 36(1), 2017. 2

[29] T. W. Sederberg and S. R. Parry. Free-form deformation of solid geometric models. In *Proceedings of the 13th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '86. ACM, 1986. 5

[30] C.-H. Shen, H. Fu, K. Chen, and S.-M. Hu. Structure recovery by part assembly. *ACM Trans. Graph.*, 31(6):180:1–180:11, Nov. 2012. 2

[31] I. Sipiran, R. Gregor, and T. Schreck. Approximate symmetry detection in partial 3d meshes. *CGF*, 33(7), 2014. 2

[32] M. Sung, V. G. Kim, R. Angst, and L. Guibas. Data-driven structural priors for shape completion. *ACM Trans. Graph.*, 34(6):175:1–175:11, 2015. 2, 7, 8

[33] S. Tulsiani, H. Su, L. J. Guibas, A. A. Efros, and J. Malik. Learning shape abstractions by assembling volumetric primitives. *CoRR*, abs/1612.00404, 2016. 8

[34] O. van Kaick, K. Xu, H. Zhang, Y. Wang, S. Sun, A. Shamir, and D. Cohen-Or. Co-hierarchical analysis of shape structures. *ACM Trans. Graph.*, 32(4):69:1–69:10, 2013. 2

[35] K. Xu, V. G. Kim, Q. Huang, N. J. Mitra, and E. Kalogerakis. Data-driven shape analysis and processing. *SIGGRAPH Asia Course notes*, 2016. 2

[36] K. Xu, H. Zhang, W. Jiang, R. Dyer, Z. Cheng, L. Liu, and B. Chen. Multi-scale partial intrinsic symmetry detection. *ACM Trans. Graph.*, 31(6):181:1–181:11, Nov. 2012. 2

[37] L. Yi, V. G. Kim, D. Ceylan, I. Shen, M. Yan, H. Su, A. Lu, Q. Huang, A. Sheffer, L. Guibas, et al. A scalable active framework for region annotation in 3d shape collections. *ACM TOG*, 35(6):210, 2016. 8

[38] L. Yi, H. Su, X. Guo, and L. Guibas. Syncspeccnn: Synchronized spectral cnn for 3d shape segmentation. 6, 8