# Learning Unsupervised Hierarchical Part Decomposition of 3D Objects from a Single RGB Image

Despoina Paschalidou[1,3,5]     Luc van Gool[3,4,5]     Andreas Geiger[1,2,5]

[1]Max Planck Institute for Intelligent Systems Tübingen

[2]University of Tübingen     [3]Computer Vision Lab, ETH Zürich     [4]KU Leuven

[5]Max Planck ETH Center for Learning Systems

{firstname.lastname}@tue.mpg.de     vangool@vision.ee.ethz.ch

## Abstract

*Humans perceive the 3D world as a set of distinct objects that are characterized by various low-level (geometry, reflectance) and high-level (connectivity, adjacency, symmetry) properties. Recent methods based on convolutional neural networks (CNNs) demonstrated impressive progress in 3D reconstruction, even when using a single 2D image as input. However, the majority of these methods focuses on recovering the local 3D geometry of an object without considering its part-based decomposition or relations between parts. We address this challenging problem by proposing a novel formulation that allows to jointly recover the geometry of a 3D object as a set of primitives as well as their latent hierarchical structure without part-level supervision. Our model recovers the higher level structural decomposition of various objects in the form of a binary tree of primitives, where simple parts are represented with fewer primitives and more complex parts are modeled with more components. Our experiments on the ShapeNet and D-FAUST datasets demonstrate that considering the organization of parts indeed facilitates reasoning about 3D geometry.*

## 1. Introduction

Within the first year of their life, humans develop a common-sense understanding of the physical behavior of the world [2]. This understanding relies heavily on the ability to properly reason about the arrangement of objects in a scene. Early works in cognitive science [22, 3, 29] stipulate that the human visual system perceives objects as a hierarchical decomposition of parts. Interestingly, while this seems to be a fairly easy task for the human brain, computer vision algorithms struggle to form such a high-level reasoning, particularly in the absence of supervision.

The structure of a scene is tightly related to the inherent hierarchical organization of its parts. At a coarse level, a
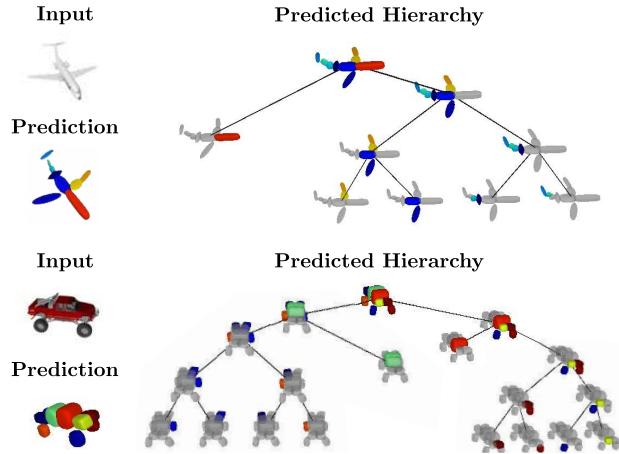


Figure 1: **Hierarchical Part Decomposition.** We consider the problem of learning structure-aware representations that go beyond part-level geometry and focus on part-level relationships. Here, we show our reconstruction as an unbalanced binary tree of primitives, given a single RGB image as input. Note that our model does not require any supervision on object parts or the hierarchical structure of the 3D object. We show that our representation is able to model different parts of an object with different levels of abstraction, leading to improved reconstruction quality.

scene can be decomposed into objects and at a finer level each object can be represented with parts and these parts with finer parts. Structure-aware representations go beyond part-level geometry and focus on global relationships between objects and object parts. In this work, we propose a structure-aware representation that considers part relationships (Fig. 1) and models object parts with multiple levels of abstraction, namely geometrically complex parts are modeled with more components and simple parts are modeled with fewer components. Such a multi-scale representation can be efficiently stored at the required level of detail, namely with less parameters (Fig. 2).

1

Recent breakthroughs in deep learning led to impressive progress in 3D shape extraction by learning a parametric function, implemented as a neural network, that maps an input image to a 3D shape represented as a mesh [34, 18, 26, 55, 60, 41], a pointcloud [14, 45, 1, 25, 51, 60], a voxel grid [5, 8, 15, 46, 47, 50, 57], 2.5D depth maps [27, 20, 44, 12] or an implicit surface [36, 7, 42, 48, 58, 37]. These approaches are mainly focused on reconstructing the geometry of an object, without taking into consideration its constituent parts. This results in non-interpretable reconstructions. To address the lack of interpretability, researchers shifted their attention to representations that employ shape primitives [53, 43, 33, 11]. While these methods yield meaningful semantic shape abstractions, part relationships do not explicitly manifest in their representations.
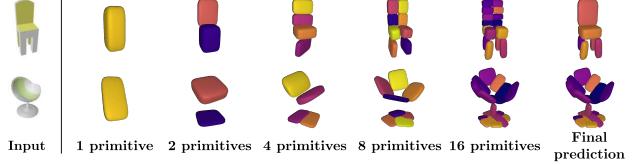
Instead of representing 3D objects as an unstructured collection of parts, we propose a novel neural network architecture that recovers the latent hierarchical layout of an object without structure supervision. In particular, we employ a neural network that learns to recursively partition an object into its constituent parts by building a latent space that encodes both the part-level hierarchy and the part geometries. The predicted hierarchical decomposition is represented as an unbalanced binary tree of primitives. More importantly, this is learned without any supervision neither on the object parts nor their structure. Instead, our model jointly infers these latent variables *during training*.

In summary, we make the following **contributions**: We jointly learn to predict part relationships and per-part geometry without any part-level supervision. The only supervision required for training our model is a watertight mesh of the 3D object. Our structure-aware representation yields semantic shape reconstructions that compare favorably to the state-of-the-art 3D reconstruction approach of [36], using significantly less parameters and without any additional post-processing. Moreover, our learned hierarchies have a semantic interpretation, as the same node in the learned tree is consistently used for representing the same object part. Experiments on the ShapeNet [6] and the Dynamic FAUST (D-FAUST) dataset [4] demonstrate the ability of our model to parse objects into structure-aware representations that are more expressive and geometrically accurate compared to approaches that only consider the 3D geometry of the object parts [53, 43, 16, 9]. Code and data is publicly available[1].

## 2. Related Work

We now discuss the most related primitive-based and structure-aware shape representations.

**Supervised Structure-Aware Representations:** Our work is related to methods that learn structure-aware shape representations that go beyond mere enumeration of object's
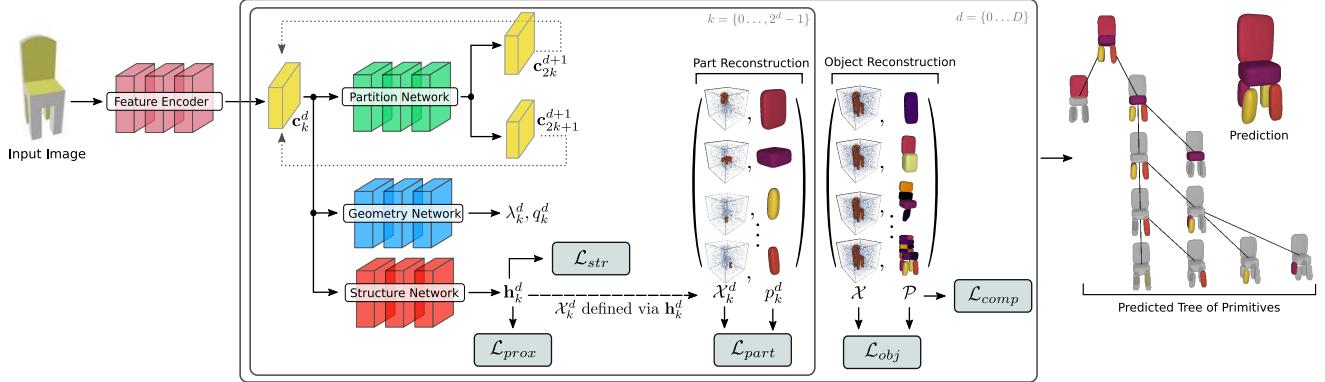
Figure 2: **Level of Detail.** Our network represents an object as a tree of primitives. At each depth level $d$, the target object is reconstructed with $2^d$ primitives, This results in a representation with various levels of detail. Naturally, reconstructions from deeper depth levels are more detailed. We associate each primitive with a unique color, thus primitives illustrated with the same color correspond to the same object part. Note that the above reconstructions are derived from the *same* model, trained with a maximum number of $2^4 = 16$ primitives. During inference, the network dynamically combines representations from different depth levels to recover the final prediction (last column).

parts and recover the higher level structural decomposition of objects based on part-level relations [38]. Li et al. [32] represent 3D shapes using a symmetry hierarchy [56] and train a recursive neural network to predict its hierarchical structure. Their network learns a hierarchical organization of bounding boxes and then fills them with voxelized parts. Note that, this model considers supervision in terms of segmentation of objects into their primitive parts. Closely related to [32] is StructureNet [39] which leverages a graph neural network to represent shapes as n-ary graphs. StructureNet considers supervision both in terms of the primitive parameters and the hierarchies. Likewise, Hu et al. [23] propose a supervised model that recovers the 3D structure of a cable-stayed bridge as a binary parsing tree. In contrast our model is unsupervised, i.e., it does not require supervision neither on the primitive parts nor the part relations.

**Physics-Based Structure-Aware Representations:** The task of inferring higher-level relationships among parts has also been investigated in different settings. Xu et al. [59] recover the object parts, their hierarchical structure and each part's dynamics by observing how objects are expected to move in the future. In particular, each part inherits the motion of its parent and the hierarchy emerges by minimizing the norm of these local displacement vectors. Kipf et al. [30] explore the use of variational autoencoders for learning the underlying interaction among various moving particles. Steenkiste et al. [54] extend the work of [17] on perceptual grouping of pixels and learn an interaction function that models whether objects interact with each other at multiple frames. For both [30, 54], the hierarchical structure emerges from interactions at multiple timestamps. In contrast to [59, 30, 54], our model does not relate hierarchies to motion, thus we do not require multiple frames for discovering the hierarchical structure.

Figure 3: **Overview.** Given an input $\mathbf{I}$ (e.g., image, voxel grid), our network predicts a binary tree of primitives $\mathcal{P}$ of maximum depth $D$. The *feature encoder* maps the input $\mathbf{I}$ into a feature vector $\mathbf{c}_0^0$. Subsequently, the *partition network* splits each feature representation $\mathbf{c}_k^d$ in two $\{\mathbf{c}_{2k}^{d+1}, \mathbf{c}_{2k+1}^{d+1}\}$, resulting in feature representations for $\{1, 2, 4, \ldots, 2^d\}$ primitives where $\mathbf{c}_k^d$ denotes the feature representation for the $k$-th primitive at depth $d$. Each $\mathbf{c}_k^d$ is passed to the *structure network* that "assigns" a part of the object to a specific primitive $p_k^d$. As a result, each $p_k^d$ is responsible for representing a specific part of the target shape, denoted as the set of points $\mathcal{X}_k^d$. Finally, the *geometry network* predicts the primitive parameters $\lambda_k^d$ and the reconstruction quality $q_k^d$ for each primitive. To compute the reconstruction loss, we measure how well the predicted primitives match the target object (*Object Reconstruction*) and the assigned parts (*Part Reconstruction*). We use plate notation to denote repetition over all nodes $k$ at each depth level $d$. The final reconstruction is shown on the right.

**Supervised Primitive-Based Representations:** Zou et al. [61] exploit LSTMs in combination with a Mixture Density Network (MDN) to learn a cuboid representation from depth maps. Similarly, Niu et al. [40] employ an RNN that iteratively predicts cuboid primitives as well as their symmetry and connectivity relationships from RGB images. More recently, Li et al. [33] utilize PointNet++ [45] for predicting per-point properties that are subsequently used for estimating the primitive parameters, by solving a series of linear least-squares problems. In contrast to [61, 40, 45], which require supervision in terms of the primitive parameters, our model is learned in an unsupervised fashion. In addition, modelling primitives with superquadrics, allows us to exploit a larger shape vocabulary that is not limited to cubes as in [61, 40] or spheres, cones, cylinders and planes as in [33]. Another line of work, complementary to ours, incorporates the principles of constructive solid geometry (CSG) [31] in a learning framework for shape modeling [49, 13, 52, 35]. These works require rich annotations for the primitive parameters and the sequence of predictions.

**Unsupervised Shape Abstraction:** Closely related to our model are the works of [53, 43] that employ a convolutional neural network (CNN) to regress the parameters of the primitives that best describe the target object, in an unsupervised manner. Primitives can be cuboids [53] or superquadrics [43] and are learned by minimizing the discrepancy between the target and the predicted shape, by either computing the truncated bi-directional distance [53] or the Chamfer-distance between points on the target and the

predicted shape [43]. While these methods learn a flat arrangement of parts, our structure-aware representation decomposes the depicted object into a hierarchical layout of semantic parts. This results in part geometries with different levels of granularity. Our model differs from [53, 43] also wrt. the optimization objective. We empirically observe that for both [53, 43], the proposed loss formulations suffer from various local minima that stem from the nature of their optimization objective. To mitigate this, we use the more robust classification loss proposed in [36, 7, 42] and train our network by learning to classify whether points lie inside or outside the target object. Very recently, [16, 9] explored such a loss for recovering shape elements from 3D objects. Genova et al. [16] leverage a CNN to learn to predict the parameters of a set of axis-aligned 3D Gaussians from a set of depth maps rendered at different viewpoints. Similarly, Deng et al. [9] employ an autoencoder to recover the geometry of an object as a collection of smooth convexes. In contrast to [16, 9], our model goes beyond the local geometry of parts and attempts to recover the underlying hierarchical structure of the object parts.

## 3. Method

In this section, we describe our novel neural network architecture for inferring structure-aware representations. Given an input $\mathbf{I}$ (e.g., RGB image, voxel grid) our goal is to learn a neural network $\phi_\theta$, which maps the input to a set of primitives that best describe the target object. The target object is represented as a set of pairs $\mathcal{X} = \{(\mathbf{x}_i, o_i)\}_{i=1}^N$,

where $\mathbf{x}_i$ corresponds to the location of the $i$-th point and $o_i$ denotes its label, namely whether $\mathbf{x}_i$ lies inside ($o_i = 1$) or outside ($o_i = 0$) the target object. We acquire these $N$ pairs by sampling points inside the bounding box of the target mesh and determine their labels using a watertight mesh of the target object. During training, our network learns to predict shapes that contain all internal points from the target mesh ($o_i = 1$) and none of the external ($o_i = 0$). We discuss our sampling strategy in our supplementary.

Instead of predicting an unstructured set of primitives, we recover a hierarchical decomposition over parts in the form of a *binary tree* of maximum depth $D$ as

$$\mathcal{P} = \{\{p_k^d\}_{k=0}^{2^d-1} \mid d = \{0 \dots D\}\} \qquad (1)$$

where $p_k^d$ is the $k$-th primitive at depth $d$. Note that for the $k$-th node at depth $d$, its parent is defined as $p_{\lfloor \frac{k}{2} \rfloor}^{d-1}$ and its two children as $p_{2k}^{d+1}$ and $p_{2k+1}^{d+1}$.

At every depth level, $\mathcal{P}$ reconstructs the target object with $\{1, 2, \dots, M\}$ primitives. $M$ is an upper limit to the maximum number of primitives and is equal to $2^D$. More specifically, $\mathcal{P}$ is constructed as follows: the root node is associated with the root primitive that represents the entire shape and is recursively split into two nodes (its children) until reaching the maximum depth $D$. This recursive partition yields reconstructions that recover the geometry of the target shape using $2^d$ primitives, where $d$ denotes the depth level (see Fig. 2). Throughout this paper, the term *node* is used interchangeably with *primitive* and always refers to the primitive associated with this particular node.

Every primitive is fully described by a set of parameters $\lambda_k^d$ that define its shape, size and position in 3D space. Since not all objects require the same number of primitives, we enable our model to predict unbalanced trees, i.e. stop recursive partitioning if the reconstruction quality is sufficient. To achieve this our network also regresses a *reconstruction quality* for each primitive denoted as $q_k^d$. Based on the value of each $q_k^d$ the network dynamically stops the recursive partitioning process resulting in parsimonious representations as illustrated in Fig. 1.

### 3.1. Network Architecture

Our network comprises three main components: (i) the *partition network* that recursively splits the shape representation into representations of parts, (ii) the *structure network* that focuses on learning the hierarchical arrangement of primitives, namely assigning parts of the object to the primitives at each depth level and (iii) the *geometry network* that recovers the primitive parameters. An overview of the proposed pipeline is illustrated in Fig. 3. The first part of our pipeline is a *feature encoder*, implemented with a ResNet-18 [21], ignoring the final fully connected layer. Instead, we only keep the feature vector of length $F = 512$ after average pooling.

**Partition Network:** The feature encoder maps the input $\mathbf{I}$ to an intermediate feature representation $\mathbf{c}_0^0 \in \mathbb{R}^F$ that describes the root node $p_0^0$. The partition network implements a function $p_\theta : \mathbb{R}^F \to \mathbb{R}^{2F}$ that recursively partitions the feature representation $\mathbf{c}_k^d$ of node $p_k^d$ into two feature representations, one for each children $\{p_{2k}^{d+1}, p_{2k+1}^{d+1}\}$:

$$p_\theta(\mathbf{c}_k^d) = \{\mathbf{c}_{2k}^{d+1}, \mathbf{c}_{2k+1}^{d+1}\}. \qquad (2)$$

Each primitive $p_k^d$ is directly predicted from $\mathbf{c}_k^d$ without considering the other intermediate features. This implies that the necessary information for predicting the primitive parameterization is entirely encapsulated in $\mathbf{c}_k^d$ and not in any other intermediate feature representation.

**Structure Network:** Due to the lack of ground-truth supervision in terms of the tree structure, we introduce the structure network that seeks to learn a pseudo-ground truth part-decomposition of the target object. More formally, it learns a function $s_\theta : \mathbb{R}^F \to \mathbb{R}^3$ that maps each feature representation $\mathbf{c}_k^d$ to $\mathbf{h}_k^d$ a spatial location in $\mathbb{R}^3$.

One can think of each $\mathbf{h}_k^d$ as the (geometric) centroid of a specific part of the target object. We define

$$\mathcal{H} = \{\{\mathbf{h}_k^d\}_{k=0}^{2^d-1} \mid d = \{0 \dots D\}\} \qquad (3)$$

the set of centroids of all parts of the object at all depth levels. From $\mathcal{H}$ and $\mathcal{X}$, we are now able to derive the part decomposition of the target object as the set of points $\mathcal{X}_k^d$ that are internal to a part with centroid $\mathbf{h}_k^d$.
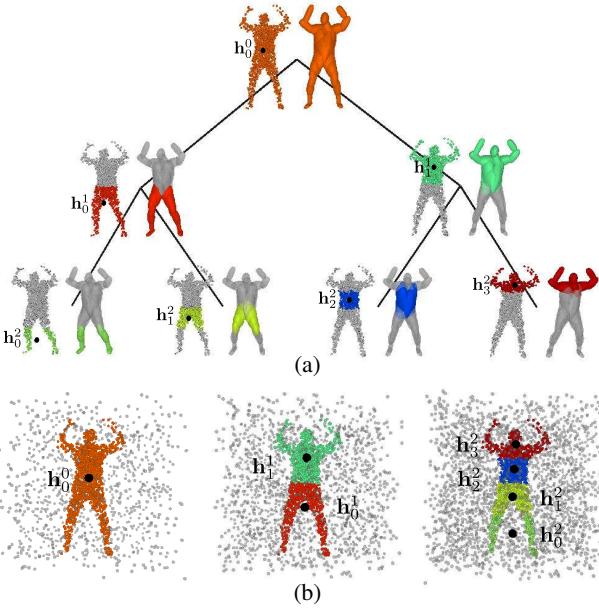
Note that, in order to learn $\mathcal{P}$, we need to be able to partition the target object into $2^d$ parts at each depth level. At the root level ($d = 0$), $\mathbf{h}_0^0$ is the centroid of the target object and $\mathcal{X}_0^0$ is equal to $\mathcal{X}$. For $d = 1$, $\mathbf{h}_0^1$ and $\mathbf{h}_1^1$ are the centroids of the two parts representing the target object. $\mathcal{X}_0^1$ and $\mathcal{X}_1^1$ comprise the same points as $\mathcal{X}_0^0$. For the external points, the labels remain the same. For the internal points, however, the labels are distributed between $\mathcal{X}_0^1$ and $\mathcal{X}_1^1$ based on whether $\mathbf{h}_0^1$ or $\mathbf{h}_1^1$ is closer. That is, $\mathcal{X}_0^1$ and $\mathcal{X}_1^1$ each contain more external labels and less internal labels compared to $\mathcal{X}_0^0$. The same process is repeated until we reach the maximum depth.

More formally, we define the set of points $\mathcal{X}_k^d$ corresponding to primitive $p_k^d$ implicitly via its centroid $\mathbf{h}_k^d$:

$$\mathcal{X}_k^d = \left\{ N_k(\mathbf{x}, o) \quad \forall (\mathbf{x}, o) \in \mathcal{X}_{\lfloor \frac{k}{2} \rfloor}^{d-1} \right\} \qquad (4)$$

Here, $\mathcal{X}_{\lfloor \frac{k}{2} \rfloor}^{d-1}$ denotes the points of the parent. The function $N_k(\mathbf{x}, o)$ assigns each $(\mathbf{x}, o) \in \mathcal{X}_{\lfloor \frac{k}{2} \rfloor}^{d-1}$ to part $p_k^d$ if it is closer to $\mathbf{h}_k^d$ than to $\mathbf{h}_{s(k)}^d$ where $s(k)$ is the sibling of $k$:

$$N_k(\mathbf{x}, o) = \begin{cases} (\mathbf{x}, 1) & \|\mathbf{h}_k^d - \mathbf{x}\| \le \|\mathbf{h}_{s(k)}^d - \mathbf{x}\| \wedge o = 1 \\ (\mathbf{x}, 0) & \text{otherwise} \end{cases}$$

$$(5)$$

Figure 4: **Structure Network.** We visualize the centroids $\mathbf{h}_k^d$ and the 3D points $\mathcal{X}_k^d$ that correspond to the estimated part $p_k^d$ for the first three levels of the tree. Fig. 4b explains visually Eq. (4). We color points based on their closest centroid $\mathbf{h}_k^d$. Points illustrated with the color associated to a part are labeled "internal" ($o = 1$). Points illustrated with gray are labeled "external" ($o = 0$).

Intuitively, this process recursively associates points to the closest sibling at each level of the binary tree where the association is determined by the label $o$. Fig. 4 illustrates the part decomposition of the target shape using $\mathcal{H}$. We visualize each part with a different color.

**Geometry Network:** The geometry network learns a function $r_\theta : \mathbb{R}^F \to \mathbb{R}^K \times [0, 1]$ that maps the feature representation $\mathbf{c}_k^d$ to its corresponding primitive parametrization $\lambda_k^d$ and the reconstruction quality prediction $q_k^d$:

$$r_\theta(\mathbf{c}_k^d) = \{\lambda_k^d, q_k^d\}. \tag{6}$$

### 3.2. Primitive Parametrization

For primitives, we use superquadric surfaces. A detailed analysis of the use of superquadrics as geometric primitives is beyond the scope of this paper, thus we refer the reader to [24, 43] for more details. Below, we focus on the properties most relevant to us. For any point $\mathbf{x} \in \mathbb{R}^3$, we can determine whether it lies inside or outside a superquadric using its implicit surface function which is commonly referred to as the *inside-outside function*:

$$f(\mathbf{x}; \lambda) = \left( \left(\frac{x}{\alpha_1}\right)^{\frac{2}{\epsilon_2}} + \left(\frac{y}{\alpha_2}\right)^{\frac{2}{\epsilon_2}} \right)^{\frac{\epsilon_2}{\epsilon_1}} + \left(\frac{z}{\alpha_3}\right)^{\frac{2}{\epsilon_1}} \tag{7}$$

where $\alpha = [\alpha_1, \alpha_2, \alpha_3]$ determine the size and $\epsilon = [\epsilon_1, \epsilon_2]$ the shape of the superquadric. If $f(\mathbf{x}; \lambda) = 1.0$, the given point $\mathbf{x}$ lies on the surface of the superquadric, if $f(\mathbf{x}; \lambda) < 1.0$ the corresponding point lies inside and if $f(\mathbf{x}; \lambda) > 1.0$ the point lies outside the superquadric. To account for numerical instabilities that arise from the exponentiations in (16), instead of directly using $f(\mathbf{x}; \lambda)$, we follow [24] and use $f(\mathbf{x}; \lambda)^{\epsilon_1}$. Finally, we convert the inside-outside function to an *occupancy function*, $g : \mathbb{R}^3 \to [0, 1]$:

$$g(\mathbf{x}; \lambda) = \sigma \left( s \left( 1 - f(\mathbf{x}; \lambda)^{\epsilon_1} \right) \right) \tag{8}$$

that results in per-point predictions suitable for the classification problem we want to solve. $\sigma(\cdot)$ is the sigmoid function and $s$ controls the sharpness of the transition of the occupancy function. To account for any rigid body motion transformations, we augment the primitive parameters with a translation vector $\mathbf{t} = [t_x, t_y, t_z]$ and a quaternion $\mathbf{q} = [q_0, q_1, q_2, q_3]$ [19], which determine the coordinate system transformation $\mathcal{T}(\mathbf{x}) = \mathbf{R}(\lambda)\mathbf{x} + \mathbf{t}(\lambda)$. Note that in (16), (17) we omit the primitive indexes $k, d$ for clarity. Visualizations of (17) are given in our supplementary.

### 3.3. Network Losses

Our optimization objective $\mathcal{L}(\mathcal{P}, \mathcal{H}; \mathcal{X})$ is a weighted sum over four terms:

$$\mathcal{L}(\mathcal{P}, \mathcal{H}; \mathcal{X}) = \mathcal{L}_{str}(\mathcal{H}; \mathcal{X}) + \mathcal{L}_{rec}(\mathcal{P}; \mathcal{X}) \\ + \mathcal{L}_{comp}(\mathcal{P}; \mathcal{X}) + \mathcal{L}_{prox}(\mathcal{P}) \tag{9}$$

**Structure Loss:** Using $\mathcal{H}$ and $\mathcal{X}$, we can decompose the target mesh into a hierarchy of disjoint parts. Namely, each $\mathbf{h}_k^d$ implicitly defines a set of points $\mathcal{X}_k^d$ that describe a specific part of the object as described in (4). To quantify how well $\mathcal{H}$ clusters the input shape $\mathcal{X}$ we minimize the sum of squared distances, similar to classical k-means:

$$\mathcal{L}_{str}(\mathcal{H}; \mathcal{X}) = \sum_{h_k^d \in \mathcal{H}} \frac{1}{2^d - 1} \sum_{(\mathbf{x}, o) \in \mathcal{X}_k^d} o \, \|\mathbf{x} - \mathbf{h}_k^d\|_2 \tag{10}$$

Note that for the loss in (10), we only consider gradients with respect to $\mathcal{H}$ as $\mathcal{X}_k^d$ is implicitly defined via $\mathcal{H}$. This results in a procedure resembling Expectation-Maximization (EM) for clustering point clouds, where computing $\mathcal{X}_k^d$ is the expectation step and each gradient updated corresponds to the maximization step. In contrast to EM, however, we minimize this loss across all instances of the training set, leading to parsimonious but consistent shape abstractions. An example of this clustering process performed at training-time is shown in Fig. 4.

**Reconstruction Loss:** The reconstruction loss measures how well the predicted primitives match the target shape. Similar to [16, 9], we formulate our reconstruction loss as a

binary classification problem, where our network learns to predict the surface boundary of the predicted shape by classifying whether points in $\mathcal{X}$ lie inside or outside the target object. To do this, we first define the occupancy function of the predicted shape at each depth level. Using the occupancy function of each primitive defined in (17), the occupancy function of the overall shape at depth $d$ becomes:

$$G^d(\mathbf{x}) = \max_{k \in 0 \ldots 2^d - 1} g_k^d \left( \mathbf{x}; \lambda_k^d \right) \qquad (11)$$

Note that (11) is simply the union of the per-primitive occupancy functions. We formulate our reconstruction loss wrt. the object and wrt. each part of the object as follows

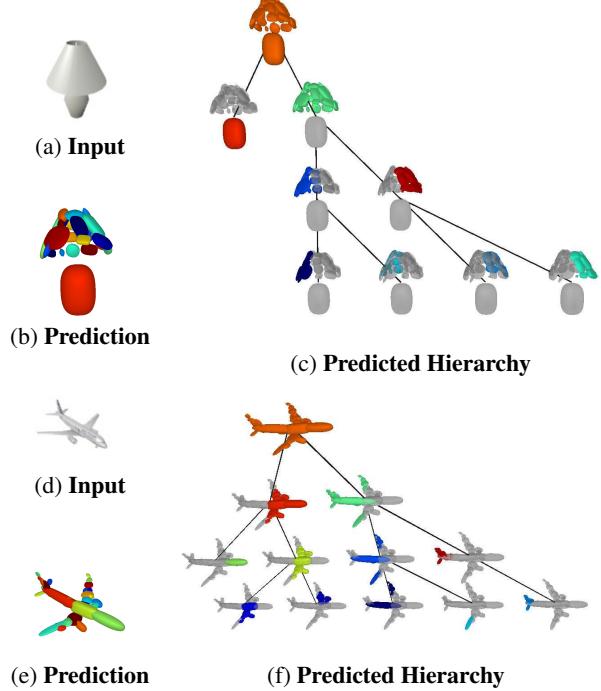$$\mathcal{L}_{rec}(\mathcal{P}; \mathcal{X}) = \sum_{(\mathbf{x},o) \in \mathcal{X}} \sum_{d=0}^{D} L \left( G^d(\mathbf{x}), o \right) + \qquad (12)$$

$$\sum_{d=0}^{D} \sum_{k=0}^{2^d - 1} \sum_{(\mathbf{x},o) \in \mathcal{X}_k^d} L \left( g_k^d \left( \mathbf{x}; \lambda_k^d \right), o \right) \quad (13)$$

where $L(\cdot)$ is the binary cross entropy loss. The first term is an *object reconstruction loss* (12) and measures how well the predicted shape at each depth level matches the target shape. The second term (13) which we refer to as *part reconstruction loss* measures how accurately each primitive $p_k^d$ matches the part of the object it represents, defined as the point set $\mathcal{X}_k^d$. Note that the *part reconstruction loss* enforces non-overlapping primitives, as $\mathcal{X}_k^d$ are non-overlapping by construction. We illustrate our reconstruction loss in Fig. 3.

**Compatibility Loss:** This loss measures how well our model is able to predict the expected *reconstruction quality* $q_k^d$ of a primitive $p_k^d$. A standard metric for measuring the reconstruction quality is the Intersection over Union (IoU). We therefore task our network to predict the *reconstruction quality* of each primitive $p_k^d$ in terms of its IoU wrt. the part of the object $\mathcal{X}_k^d$ it represents:

$$\mathcal{L}_{comp}(\mathcal{P}; \mathcal{X}) = \sum_{d=0}^{\mathcal{D}} \sum_{k=0}^{2^d - 1} \left( q_k^d - \text{IoU}(p_k^d, \mathcal{X}_k^d) \right)^2 \qquad (14)$$

During inference, $q_k^d$ allows for further partitioning primitives whose IoU is below a threshold $q_{th}$ and to stop if the reconstruction quality is high (the primitive fits the object part well). As a result, our model predicts an unbalanced tree of primitives where objects can be represented with various number of primitives from 1 to $2^D$. This results in parsimonious representations where simple parts are represented with fewer primitives. We empirically observe that the threshold value $q_{th}$ does not significantly affect our results, thus we empirically set it to 0.6. During training, we do not use the predicted reconstruction quality $q_k^d$ to dynamically partition the nodes but instead predict the full tree.



(a) **Input**

(b) **Prediction**

(c) **Predicted Hierarchy**

(d) **Input**

(e) **Prediction**

(f) **Predicted Hierarchy**

Figure 5: **Predicted Hierarchies on ShapeNet**. Our model recovers the geometry of an object as an unbalanced hierarchy over primitives, where simpler parts (e.g. base of the lamp) are represented with few primitives and more complex parts (e.g. wings of the plane) with more.

**Proximity Loss:** This term is added to counteract vanishing gradients due to the sigmoid in (17). For example, if the initial prediction of a primitive is far away from the target object, the reconstruction loss will be large while its gradients will be small. As a result, it is impossible to "move" this primitive to the right location. Thus, we introduce a proximity loss which encourages the center of each primitive $p_k^d$ to be close to the centroid of the part it represents:

$$\mathcal{L}_{prox}(\mathcal{P}) = \sum_{d=0}^{D} \sum_{k=0}^{2^d - 1} \| \mathbf{t}(\lambda_k^d) - \mathbf{h}_k^d \|_2 \qquad (15)$$

where $\mathbf{t}(\lambda_k^d)$ is the translation vector of the primitive $p_k^d$ and $\mathbf{h}_k^d$ is the centroid of the part it represents. We demonstrate the vanishing gradient problem in our supplementary.

## 4. Experiments

In this section, we provide evidence that our structure-aware representation yields semantic shape abstractions while achieving competitive (or even better results) than various state-of-the-art shape reconstruction methods, such as [36]. Moreover, we also investigate the quality of the learned hierarchies and show that the use of our structure-aware representation yields semantic scene parsings. Im-

| | Chamfer-$L_1$ | | | | | IoU | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Category | OccNet [36] | SQs [43] | SIF [16] | CvxNets [9] | Ours | OccNet [36] | SQs [43] | SIF [16] | CvxNets [9] | Ours |
| airplane | 0.147 | 0.122 | **0.065** | 0.093 | 0.175 | 0.571 | 0.456 | 0.530 | **0.598** | 0.529 |
| bench | 0.155 | **0.114** | 0.131 | 0.133 | 0.153 | 0.485 | 0.202 | 0.333 | **0.461** | 0.437 |
| cabinet | 0.167 | **0.087** | 0.102 | 0.160 | **0.087** | 0.733 | 0.110 | 0.648 | **0.709** | 0.658 |
| car | 0.159 | 0.117 | **0.056** | 0.103 | 0.141 | 0.737 | 0.650 | 0.657 | 0.675 | **0.702** |
| chair | 0.228 | 0.138 | 0.192 | 0.337 | **0.114** | 0.501 | 0.176 | 0.389 | 0.491 | **0.526** |
| display | 0.278 | **0.106** | 0.208 | 0.223 | 0.137 | 0.471 | 0.200 | 0.491 | 0.576 | **0.633** |
| lamp | 0.479 | 0.189 | 0.454 | 0.795 | **0.169** | 0.371 | 0.189 | 0.260 | 0.311 | **0.441** |
| speaker | 0.300 | 0.132 | 0.253 | 0.462 | **0.108** | 0.647 | 0.136 | 0.577 | 0.620 | **0.660** |
| rifle | 0.141 | 0.127 | **0.069** | 0.106 | 0.203 | 0.474 | **0.519** | 0.463 | 0.515 | 0.435 |
| sofa | 0.194 | **0.106** | 0.146 | 0.164 | 0.128 | 0.680 | 0.122 | 0.606 | 0.677 | **0.693** |
| table | 0.189 | **0.110** | 0.264 | 0.358 | 0.122 | 0.506 | 0.180 | 0.372 | 0.473 | **0.491** |
| phone | 0.140 | 0.112 | **0.095** | 0.083 | 0.149 | 0.720 | 0.185 | 0.658 | 0.719 | **0.770** |
| vessel | 0.218 | 0.125 | **0.108** | 0.173 | 0.178 | 0.530 | 0.471 | 0.502 | 0.552 | **0.570** |
| mean | 0.215 | **0.122** | 0.165 | 0.245 | 0.143 | 0.571 | 0.277 | 0.499 | 0.567 | **0.580** |

Table 1: **Single Image Reconstruction on ShapeNet.** Quantitative evaluation of our method against OccNet [36] and primitive-based methods with superquadrics [43] (SQs), SIF [16] and CvxNets [9]. We report the volumeteric IoU (higher is better) and the Chamfer-$L_1$ distance (lower is better) wrt. the ground-truth mesh.
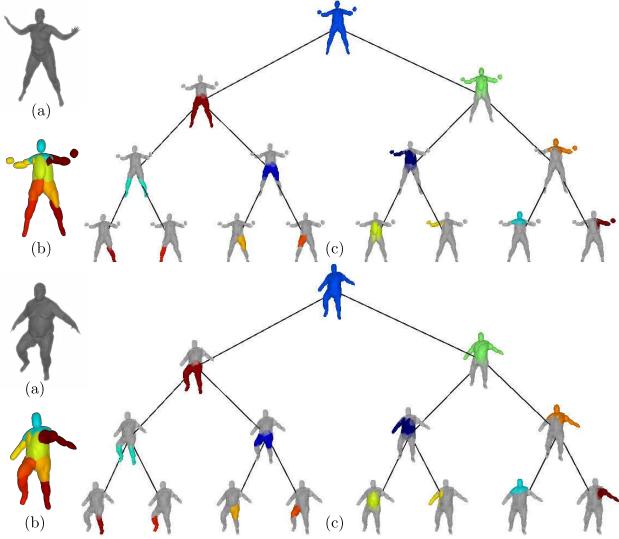


Figure 6: **Predicted Hierarchies on D-FAUST.** We visualize the input RGB image (a), the prediction (b) and the predicted hierarchy (c). We associate each primitive with a color and we observe that our network learns semantic mappings of body parts across different articulations, e.g. node $(3, 3)$ is used for representing the upper part of the left leg, whereas node $(1, 1)$ is used for representing the upper body.

plementation details and ablations on the impact of various components of our model are detailed in the supplementary.

**Datasets:** First, we use the ShapeNet [6] subset of Choy et al. [8], training our model using the same image renderings and train/test splits as Choy et al. Furthermore, we also

experiment with the Dynamic FAUST (D-FAUST) dataset [4], which contains meshes for 129 sequences of 10 humans performing various tasks, such as "running", "punching" or "shake arms". We randomly divide these sequences into training (91), test (29) and validation (9).

**Baselines:** Closely related to our work are the shape parsing methods of [53] and [43] that employ cuboids and superquadric surfaces as primitives. We refer to [43] as SQs and we evaluate using their publicly available code[2]. Moreover, we also compare to the Structured Implicit Function (SIF) [16] that represent the object's geometry as the isolevel of the sum of a set of Gaussians and to the CvxNets [9] that represent the object parts using smooth convex shapes. Finally, we also report results for OccNet [36], which is the state-of-the-art implicit shape reconstruction technique. Note that in contrast to us, [36] does not consider part decomposition or any form of latent structure.

**Evaluation Metrics:** Similar to [36, 16, 9], we evaluate our model quantitatively and report the mean Volumetric IoU and the Chamfer-$L_1$ distance. Both metrics are discussed in detail in our supplementary.

## 4.1. Results on ShapeNet

We evaluate our model on the single-view 3D reconstruction task and compare against various state-of-the-art methods. We follow the standard experimental setup and train a single model for the 13 ShapeNet objects. Both our model and [43] are trained for a maximum number of 64

---

[2]https://superquadrics.com

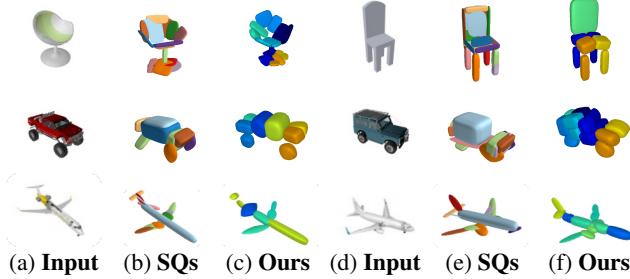(a) **Input** (b) **SQs** (c) **Ours** (d) **Input** (e) **SQs** (f) **Ours**

Figure 7: **Single Image 3D Reconstruction.** The input image is shown in (a, d), the other columns show the results of our method (c, f) compared to [43] (b, e). Additional qualitative results are provided in the supplementary.

primitives ($D = 6$). For SIF [16] and CvxNets [9] the reported results are computed using 50 shape elements. The quantitative results are reported in Table 1. We observe that our model outperforms the primitive-based baselines in terms of the IoU as well as the OccNet [36] for the majority of objects ($7/13$). Regarding Chamfer-$L_1$, our model is the second best amongst primitive representations, as [43] is optimized for this metric. This also justifies that [43] performs worse in terms of IoU. While our model performs on par with existing state-of-the-art primitive representations in terms of Chamfer-$L_1$, it also recovers hierarchies, which none of our baselines do. A qualitative comparison of our model with SQs [43] is depicted in Fig. 7. Fig. 5 visualizes the learned hierarchy for this model. We observe that our model recovers unbalanced binary trees that decompose a 3D object into a set of parts. Note that [53, 43] were originally introduced for volumetric 3D reconstruction, thus we provide an experiment on this task in our supplementary.

### 4.2. Results on D-FAUST

We also demonstrate results on the Dynamic FAUST (D-FAUST) dataset [4], which is very challenging due to the fine structure of the human body. We evaluate our model on the single-view 3D reconstruction task and compare with [43]. Both methods are trained for a maximum number of 32 primitives ($D = 5$). Fig. 6 illustrates the predicted hierarchy on different humans from the test set. We note that the predicted hierarchies are indeed semantic, as the same nodes are used for modelling the same part of the human body. Fig. 8 compares the predictions of our model with SQs. We observe that while our baseline yields more parsimonious abstractions, their level of detail is limited. On the contrary, our model captures the geometry of the human body with more detail. This is also validated quantitatively, from Table 2. Note that in contrast to ShapeNet, D-FAUST does not contain long, thin (e.g. legs of tables, chairs) or hollow parts (e.g. cars), thus optimizing for either Chamfer-L1 or IoU leads to similar results. Hence, our method out-



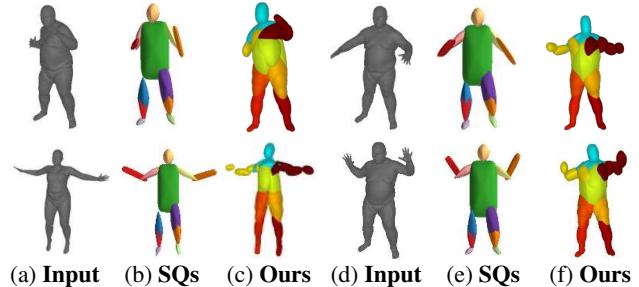(a) **Input** (b) **SQs** (c) **Ours** (d) **Input** (e) **SQs** (f) **Ours**

Figure 8: **Single Image 3D Reconstruction**. Qualitative comparison of our reconstructions (c, f), to [43] that does not consider any form of structure (b, e). The input RGB image is shown in (a, d). Note how our representation yields geometrically more accurate reconstructions, while being semantic, e.g., the primitive colored in blue consistently represents the head of the human while the primitive colored in orange captures the left thigh. Additional qualitative results are provided in the supplementary.

|  | IoU | Chamfer-$L_1$ |
|---|---|---|
| SQs [43] | 0.608 | 0.189 |
| Ours | **0.699** | **0.098** |

Table 2: **Single Image Reconstruction on D-FAUST.** We report the volumetric IoU and the Chamfer-L1 wrt. the ground-truth mesh for our model compared to [43].

performs [43] also in terms of Chamfer-L1. Due to lack of space, we only illustrate the predicted hierarchies up to the fourth depth level. The full hierarchies are provided in the supplementary.

## 5. Conclusion

We propose a learning-based approach that jointly predicts part relationships together with per-part geometries in the form of a binary tree without requiring any part-level annotations for training. Our model yields geometrically accurate primitive-based reconstructions that outperform existing shape abstraction techniques while performing competitively with more flexible implicit shape representations. In future work, we plan to to extend our model and predict hierarchical structures that remain consistent in time, thus yielding kinematic trees of objects. Another future direction, is to consider more flexible primitives such as general convex shapes and incorporate additional constraints e.g. symmetry to further improve the reconstructions.

## Acknowledgments

# References

[1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J. Guibas. Learning representations and generative models for 3d point clouds. In *Proc. of the International Conf. on Machine learning (ICML)*, 2018. 2

[2] Renée Baillargeon. Infants' physical world. *Current directions in psychological science*, 13(3):89–94, 2004. 1

[3] Irving Biederman. Human image understanding: Recent research and a theory. *Computer Vision, Graphics, and Image Processing*, 1986. 1

[4] Federica Bogo, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Dynamic FAUST: registering human bodies in motion. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2, 7, 8, 12, 22

[5] André Brock, Theodore Lim, James M. Ritchie, and Nick Weston. Generative and discriminative voxel modeling with convolutional neural networks. *arXiv.org*, 1608.04236, 2016. 2

[6] Angel X. Chang, Thomas A. Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qi-Xing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository. *arXiv.org*, 1512.03012, 2015. 2, 7, 12, 21

[7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 3

[8] Christopher Bongsoo Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2016. 2, 7, 21

[9] Boyang Deng, Kyle Genova, Soroosh Yazdani, Sofien Bouaziz, Geoffrey Hinton, and Andrea Tagliasacchi. Cvxnets: Learnable convex decomposition. *arXiv.org*, 2019. 2, 3, 5, 7, 8, 21

[10] Jia Deng, Wei Dong, Richard Socher, Li jia Li, Kai Li, and Li Fei-fei. Imagenet: A large-scale hierarchical image database. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2009. 16

[11] Theo Deprelle, Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. Learning elementary structures for 3d shape generation and matching. In *Advances in Neural Information Processing Systems (NIPS)*, 2019. 2

[12] Simon Donne and Andreas Geiger. Learning non-volumetric depth fusion using successive reprojections. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[13] Kevin Ellis, Daniel Ritchie, Armando Solar-Lezama, and Joshua B. Tenenbaum. Learning to infer graphics programs from hand-drawn images. In *Advances in Neural Information Processing Systems (NIPS)*, 2018. 3

[14] Haoqiang Fan, Hao Su, and Leonidas J. Guibas. A point set generation network for 3d object reconstruction from a single image. *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2

[15] Matheus Gadelha, Subhransu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. In *Proc. of the International Conf. on 3D Vision (3DV)*, 2017. 2

[16] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T. Freeman, and Thomas A. Funkhouser. Learning shape templates with structured implicit functions. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2, 3, 5, 7, 8, 21

[17] Klaus Greff, Sjoerd van Steenkiste, and Jürgen Schmidhuber. Neural expectation maximization. In *Advances in Neural Information Processing Systems (NIPS)*, 2017. 2

[18] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, and Mathieu Aubry. AtlasNet: A papier-mâché approach to learning 3d surface generation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[19] William Rowan Hamilton. Xi. on quaternions; or on a new system of imaginaries in algebra. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 33(219):58–60, 1848. 5

[20] Wilfried Hartmann, Silvano Galliani, Michal Havlena, Luc Van Gool, and Konrad Schindler. Learned multi-patch similarity. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2017. 2

[21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016. 4, 16

[22] Donald D Hoffman and Whitman A Richards. Parts of recognition. *Cognition*, 18(1-3):65–96, 1984. 1

[23] Fangqiao Hu, Jin Zhao, Yong Hunag, and Hui Li. Learning structural graph layouts and 3d shapes for long span bridges 3d reconstruction. *arXiv.org*, abs/1907.03387, 2019. 2

[24] Ales Jaklic, Ales Leonardis, and Franc Solina. *Segmentation and Recovery of Superquadrics*, volume 20 of *Computational Imaging and Vision*. Springer, 2000. 5, 12

[25] Li Jiang, Shaoshuai Shi, Xiaojuan Qi, and Jiaya Jia. GAL: geometric adversarial loss for single-view 3d-object reconstruction. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2018. 2

[26] Angjoo Kanazawa, Shubham Tulsiani, Alexei A. Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2018. 2

[27] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. In *Advances in Neural Information Processing Systems (NIPS)*, 2017. 2

[28] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proc. of the International Conf. on Learning Representations (ICLR)*, 2015. 17

[29] Katherine D Kinzler and Elizabeth S Spelke. Core systems in human cognition. *Progress in brain research*, 164:257–264, 2007. 1

[30] Thomas Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling, and Richard Zemel. Neural relational inference for interacting systems. In *Proc. of the International Conf. on Machine learning (ICML)*, 2018. 2

[31] David H Laidlaw, W Benjamin Trumbore, and John F Hughes. Constructive solid geometry for polyhedral objects. In *ACM Trans. on Graphics*, 1986. 3

[32] Jun Li, Kai Xu, Siddhartha Chaudhuri, Ersin Yumer, Hao (Richard) Zhang, and Leonidas J. Guibas. GRASS: generative recursive autoencoders for shape structures. *ACM Trans. on Graphics*, 36(4), 2017. 2

[33] Lingxiao Li, Minhyuk Sung, Anastasia Dubrovina, Li Yi, and Leonidas Guibas. Supervised fitting of geometric primitives to 3d point clouds. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 3

[34] Yiyi Liao, Simon Donne, and Andreas Geiger. Deep marching cubes: Learning explicit surface representations. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[35] Yunchao Liu, Zheng Wu, Daniel Ritchie, William T Freeman, Joshua B Tenenbaum, and Jiajun Wu. Learning to describe scenes with programs. In *Proc. of the International Conf. on Learning Representations (ICLR)*, 2019. 3

[36] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 3, 6, 7, 8, 18

[37] Mateusz Michalkiewicz, Jhony Kaesemodel Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders P. Eriksson. Implicit surface representations as layers in neural networks. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2

[38] Niloy J. Mitra, Michael Wand, Hao Zhang, Daniel Cohen-Or, Vladimir G. Kim, and Qi-Xing Huang. Structure-aware shape processing. In *ACM Trans. on Graphics*, pages 13:1–13:21, 2014. 2

[39] Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas Guibas. Structurenet: Hierarchical graph networks for 3d shape generation. In *ACM Trans. on Graphics*, 2019. 2

[40] Chengjie Niu, Jun Li, and Kai Xu. Im2struct: Recovering 3d shape structure from a single RGB image. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3

[41] Junyi Pan, Xiaoguang Han, Weikai Chen, Jiapeng Tang, and Kui Jia. Deep mesh reconstruction from single RGB images via topology modification networks. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2

[42] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 3

[43] Despoina Paschalidou, Ali Osman Ulusoy, and Andreas Geiger. Superquadrics revisited: Learning 3d shape parsing beyond cuboids. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 3, 5, 7, 8, 16, 18, 21, 22, 23, 24, 25, 27

[44] Despoina Paschalidou, Ali Osman Ulusoy, Carolin Schmitt, Luc van Gool, and Andreas Geiger. Raynet: Learning volumetric 3d reconstruction with ray potentials. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[45] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems (NIPS)*, 2017. 2, 3

[46] Danilo Jimenez Rezende, S. M. Ali Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, and Nicolas Heess. Unsupervised learning of 3d structure from images. In *Advances in Neural Information Processing Systems (NIPS)*, 2016. 2

[47] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2

[48] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2

[49] Gopal Sharma, Rishabh Goyal, Difan Liu, Evangelos Kalogerakis, and Subhransu Maji. Csgnet: Neural shape parser for constructive solid geometry. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3

[50] David Stutz and Andreas Geiger. Learning 3d shape completion from laser scan data with weak supervision. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[51] Hugues Thomas, Charles R. Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J. Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2

[52] Yonglong Tian, Andrew Luo, Xingyuan Sun, Kevin Ellis, William T Freeman, Joshua B Tenenbaum, and Jiajun Wu. Learning to infer and execute 3d shape programs. In *Proc. of the International Conf. on Learning Representations (ICLR)*, 2019. 3

[53] Shubham Tulsiani, Hao Su, Leonidas J. Guibas, Alexei A. Efros, and Jitendra Malik. Learning shape abstractions by assembling volumetric primitives. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2, 3, 7, 8, 18, 21, 22

[54] Sjoerd van Steenkiste, Michael Chang, Klaus Greff, and Jürgen Schmidhuber. Relational neural expectation maximization: Unsupervised discovery of objects and their interactions. In *Proc. of the International Conf. on Learning Representations (ICLR)*, 2018. 2

[55] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2018. 2

[56] Yanzhen Wang, Kai Xu, Jun Li, Hao Zhang, Ariel Shamir, Ligang Liu, Zhi-Quan Cheng, and Yueshan Xiong. Symmetry hierarchy of man-made objects. In *EUROGRAPHICS*, 2011. 2

[57] Haozhe Xie, Hongxun Yao, Xiaoshuai Sun, Shangchen Zhou, and Shengping Zhang. Pix2vox: Context-aware 3d reconstruction from single and multi-view images. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2

[58] Qiangeng Xu, Weiyue Wang, Duygu Ceylan, Radomír Mech, and Ulrich Neumann. DISN: deep implicit surface network for high-quality single-view 3d reconstruction. In *Advances in Neural Information Processing Systems (NIPS)*, pages 490–500, 2019. 2

[59] Zhenjia Xu, Zhijian Liu, Chen Sun, Kevin Murphy, William T Freeman, Joshua B Tenenbaum, and Jiajun Wu. Unsupervised discovery of parts, structure, and dynamics. In *Proc. of the International Conf. on Learning Representations (ICLR)*, 2019. 2

[60] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge J. Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2

[61] C. Zou, E. Yumer, J. Yang, D. Ceylan, and D. Hoiem. 3d-prnn: Generating shape primitives with recurrent neural networks. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2017. 3

# Supplementary Material for
# Learning Unsupervised Hierarchical Part Decomposition of 3D Objects from a Single RGB Image

Despoina Paschalidou[1,3,5]    Luc van Gool[3,4,5]    Andreas Geiger[1,2,5]

[1]Max Planck Institute for Intelligent Systems Tübingen
[2]University of Tübingen    [3]Computer Vision Lab, ETH Zürich    [4]KU Leuven
[5]Max Planck ETH Center for Learning Systems

{firstname.lastname}@tue.mpg.de    vangool@vision.ee.ethz.ch

## Abstract

*In this **supplementary document**, we first present examples of our occupancy function. In addition, we present a detailed overview of our network architecture and the training procedure. We then discuss how various components influence the performance of our model on the single-view 3D reconstruction task. Finally, we provide additional experimental results on more categories from the ShapeNet dataset [6] and on the D-FAUST dataset [4] together with the corresponding hierarchical structures. The **supplementary video** shows 3D animations of the predicted structural hierarchy for various objects from the ShapeNet dataset as well as humans from the D-FAUST.*

## A. Occupancy Function

In this section, we provide illustrations of the occupancy function $g$ for different primitive parameters and for different sharpness values. For any point $\mathbf{x} \in \mathbb{R}^3$, we can determine whether it lies inside or outside a superquadric using its implicit surface function which is commonly referred to as the *inside-outside function*:

$$f(\mathbf{x}; \lambda) = \left( \left( \frac{x}{\alpha_1} \right)^{\frac{2}{\epsilon_2}} + \left( \frac{y}{\alpha_2} \right)^{\frac{2}{\epsilon_2}} \right)^{\frac{\epsilon_2}{\epsilon_1}} + \left( \frac{z}{\alpha_3} \right)^{\frac{2}{\epsilon_1}} \tag{16}$$

where $\alpha = [\alpha_1, \alpha_2, \alpha_3]$ determine the size and $\epsilon = [\epsilon_1, \epsilon_2]$ determine the shape of the superquadric. If $f(\mathbf{x}; \lambda) = 1.0$, the given point $\mathbf{x}$ lies on the surface of the superquadric, if $f(\mathbf{x}; \lambda) < 1.0$ the corresponding point lies inside and if $f(\mathbf{x}; \lambda) > 1.0$ the point lies outside the superquadric. To account for numerical instabilities that arise from the exponentiations in (16), instead of directly using $f(\mathbf{x}; \lambda)$, we follow [24] and use $f(\mathbf{x}; \lambda)^{\epsilon_1}$. In addition, we also convert the inside-outside function
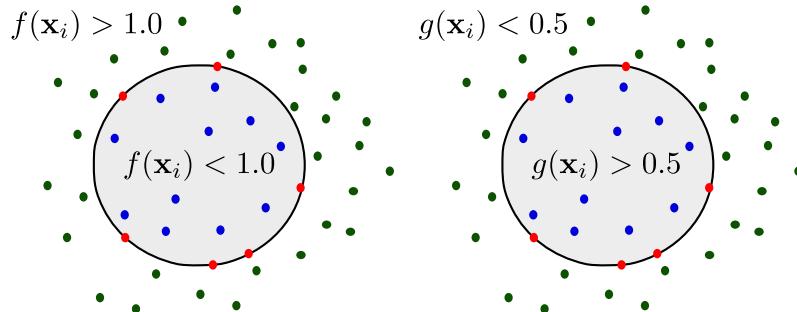


Figure 9: **Implicit surface function of superquadrics.** We visualize the 2D slice of $f(\mathbf{x}_i)$ and $g(\mathbf{x}_i)$ for a superquadric with $\alpha_1 = \alpha_2 = \alpha_3 = \epsilon_1 = \epsilon_2 = 1$.

(a) $\epsilon_1 = 0.25$, and $\epsilon_2 = 0.25$

(b) $\epsilon_1 = 0.25$, and $\epsilon_2 = 0.5$

(c) $\epsilon_1 = 0.25$, and $\epsilon_2 = 1.0$

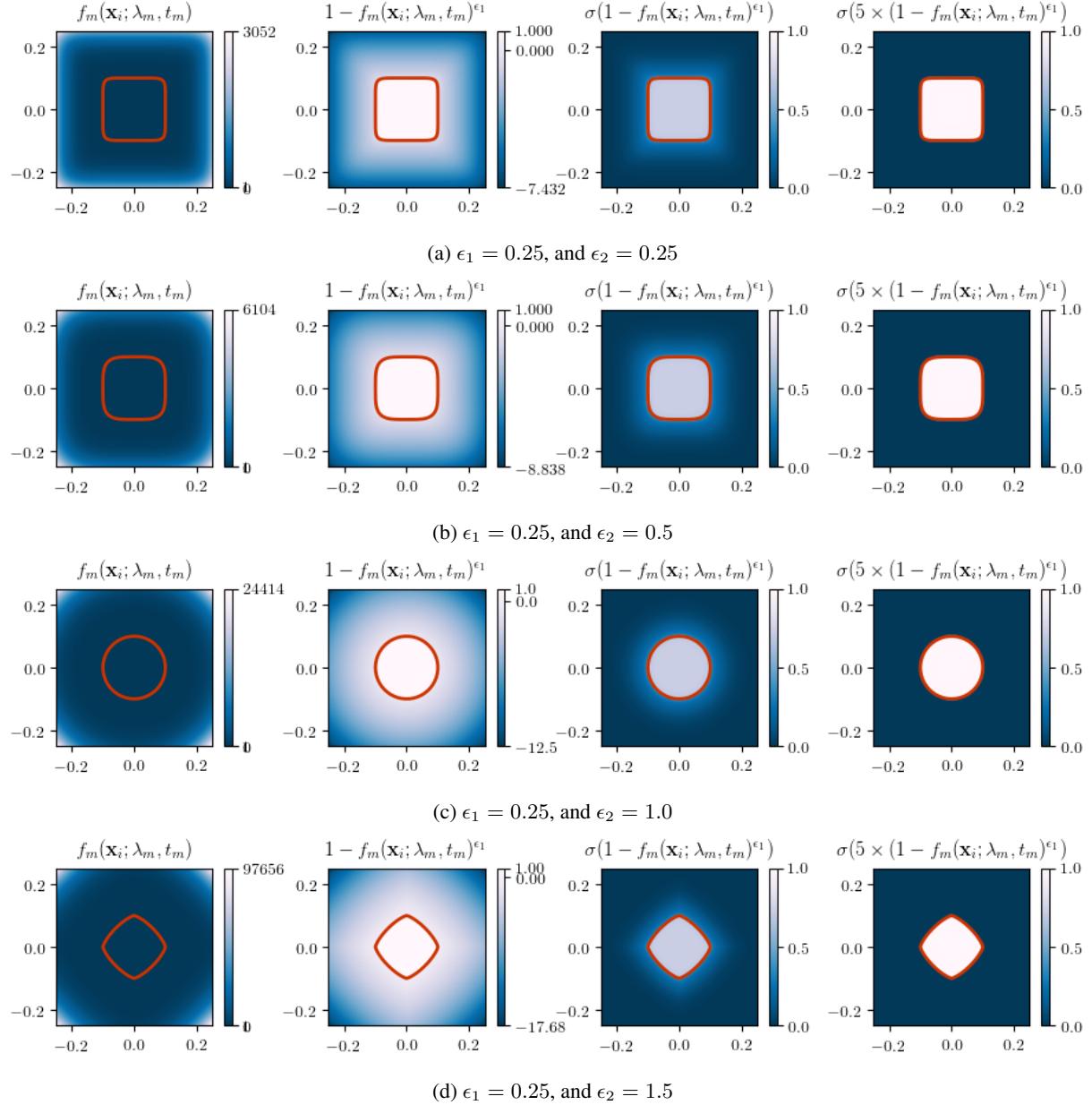(d) $\epsilon_1 = 0.25$, and $\epsilon_2 = 1.5$

Figure 10: **Implicit surface function** We visualize the implicit surface function for different primitive parameters and for different sharpness values. The surface boundary is drawn with red.

to an *occupancy function*, $g : \mathbb{R}^3 \to [0, 1]$:

$$g(\mathbf{x}; \lambda) = \sigma \left( s \left( 1 - f(\mathbf{x}; \lambda)^{\epsilon_1} \right) \right) \tag{17}$$

that results in per-point predictions suitable for the classification problem we want to solve. $\sigma(\cdot)$ is the sigmoid function and $s$ controls the sharpness of the transition of the occupancy function. As a result, if $g(\mathbf{x}; \lambda) < 0.5$ the corresponding point lies outside and if $g(\mathbf{x}; \lambda) > 0.5$ the point lies inside the superquadric. Fig. 9 visualizes the range of the implicit surface function of superquadrics of (16) and (17). Fig. 10+11+12 visualize the implicit surface function for different values of $\epsilon_1$ and $\epsilon_2$ and different values of sharpness $s$. We observe that without applying the sigmoid to (16) the range of values of (16) varies significantly for different primitive parameters.
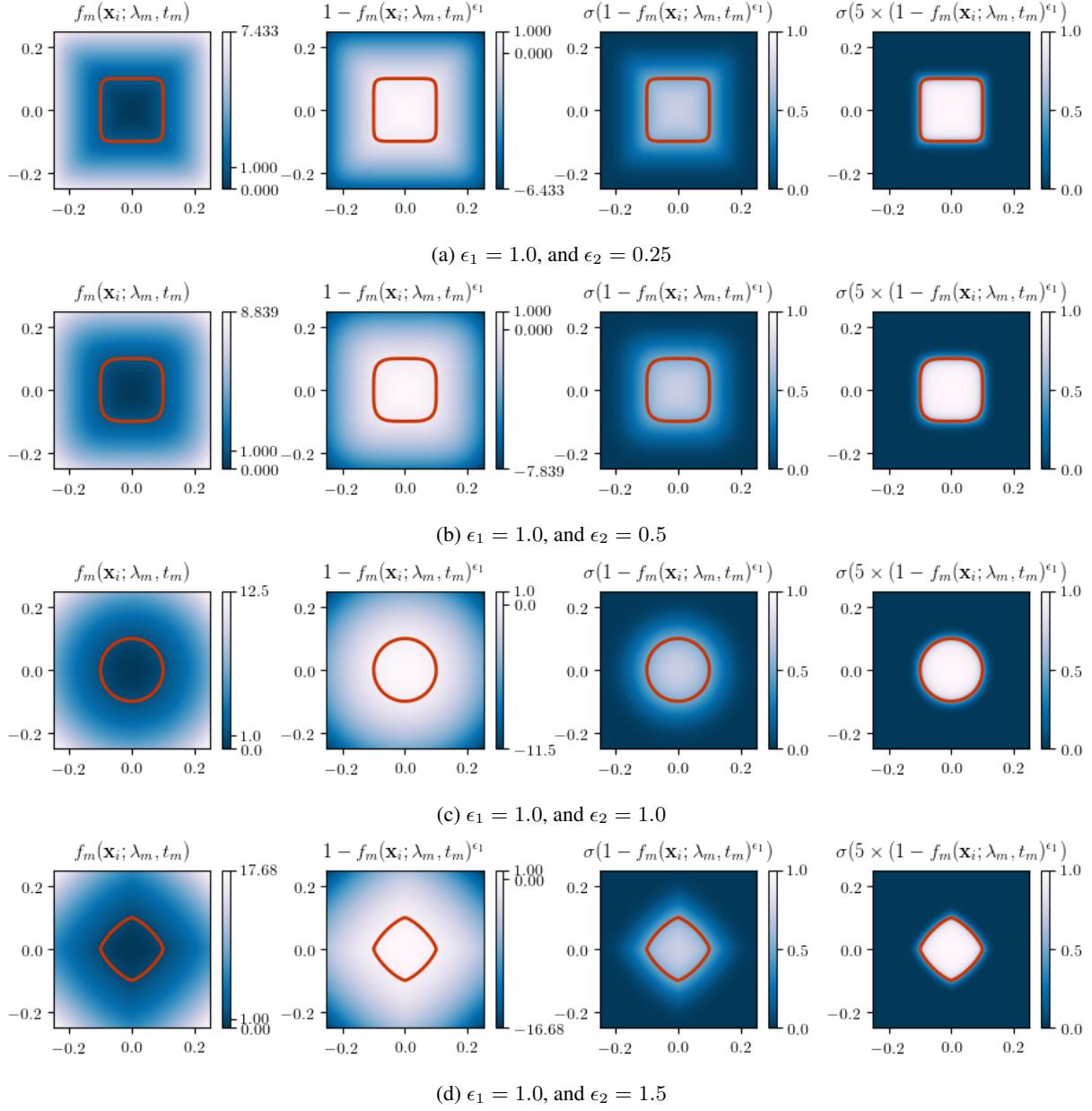
Figure 11: **Implicit surface function** We visualize the implicit surface function for different primitive parameters and for different sharpness values. The surface boundary is drawn with red.
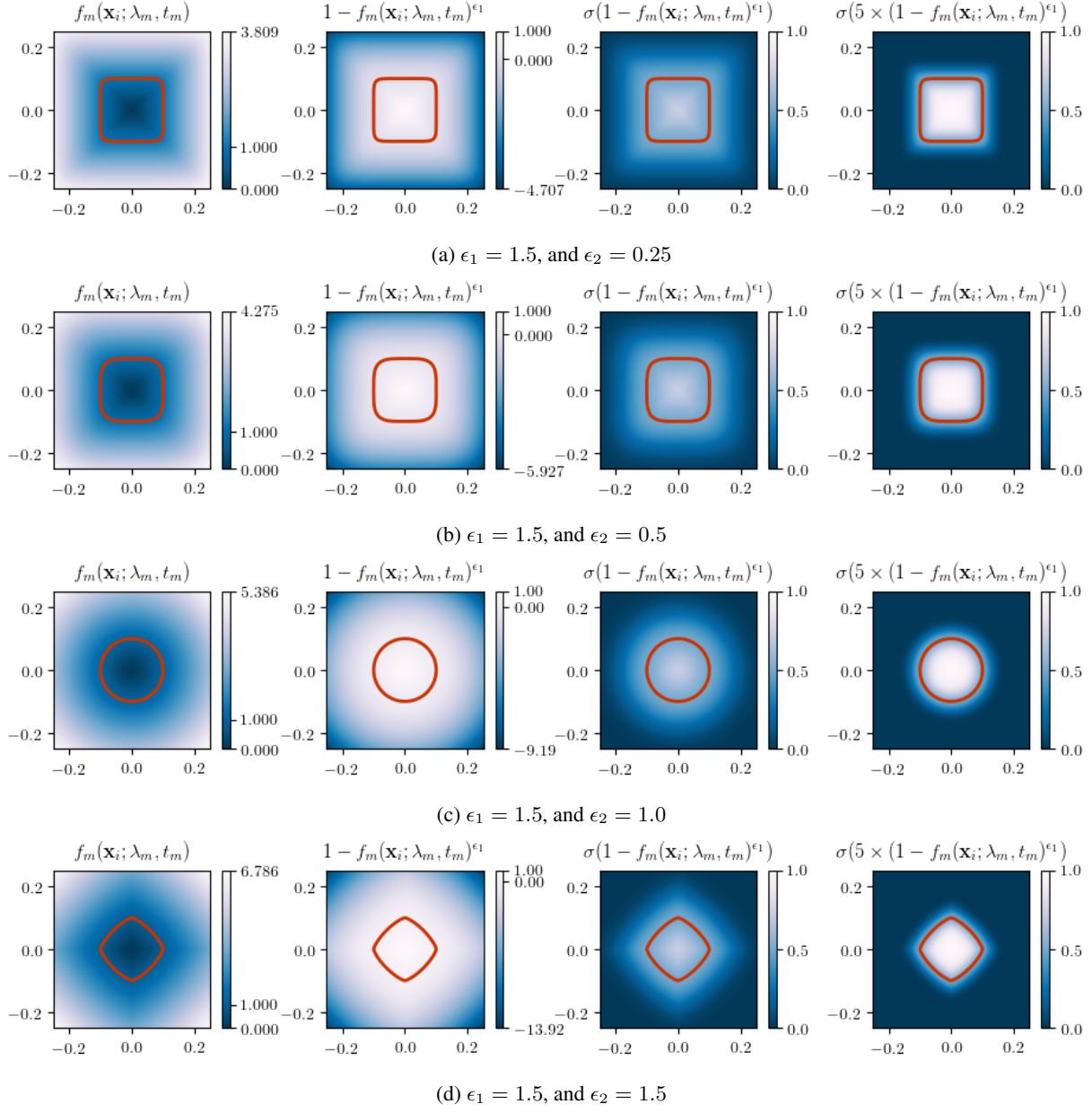
Figure 12: **Implicit surface function** We visualize the implicit surface function for different primitive parameters and for different sharpness values. The surface boundary is drawn with red.

## B. Implementation Details

In this section, we provide a detailed description of our network architecture. We then describe our sampling strategy and provide details on the metrics we use both for training and testing. Finally, we show how various components influence the performance of our model on the single-view 3D reconstruction task.

### B.1. Network Architecture

Here we describe the architecture of each individual component of our model, shown in Figure 3 of our main submission.

**Feature Encoder:** The feature encoder depends on the type of the input, namely whether it is an image or a binary occupancy grid. For the single view 3D reconstruction task, we use a ResNet-18 architecture [21] (Fig. 13a), which was pretrained on the ImageNet dataset [10]. From the original design, we ignore the final fully connected layer keeping only the feature vector of length $F = 512$ after average pooling. For the volumetric 3D reconstruction task, where the input is a binary occupancy grid, we use the feature encoder proposed in [43](Fig. 13b). Note that the feature encoder is used as a generic feature extractor from the input representation.



(a) Single-view 3D Reconstruction
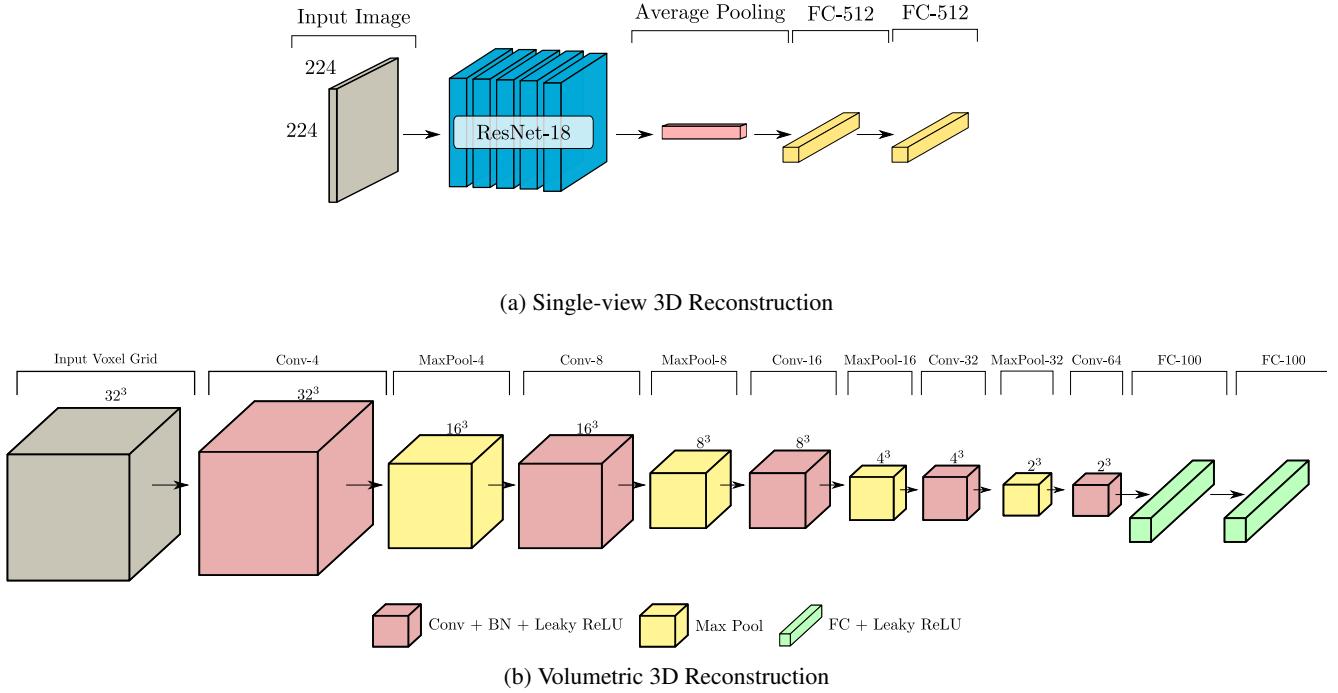


(b) Volumetric 3D Reconstruction

Figure 13: **Feature Encoder Architectures.** Depending on the type of the input, we employ two different network architectures. (a) For the single view 3D reconstruction task we use a ResNet-18 architecture [21] (b) For a binary occupancy grid as an input, we leverage the network architecture of [43].

**Partition Network:** The partition network implements a function $p_\theta : \mathbb{R}^F \to \mathbb{R}^{2F}$ that recursively partitions the feature representation $\mathbf{c}_k^d$ of node $p_k^d$ into two feature representations, one for each child $\{p_{2k}^{d+1}, p_{2k+1}^{d+1}\}$. The partition network (Fig. 14a) comprises two fully connected layers followed by RELU non linearity.

**Structure Network:** The structure network maps each feature representation $\mathbf{c}_k^d$ to $\mathbf{h}_k^d$ a spatial location in $\mathbb{R}^3$. The structure network (Fig. 14b) consists of two fully connected layers followed by RELU non linearity.

**Geometry Network:** The geometry network learns a function $r_\theta : \mathbb{R}^F \to \mathbb{R}^K \times [0, 1]$ that maps the feature representation $\mathbf{c}_k^d$ to its corresponding primitive parametrization $\lambda_k^d$ and the reconstruction quality prediction $q_k^d$. In particular, the geometry network consists of five regressors that predict the parameters of the superquadrics (size $\alpha$, shape $\epsilon$ and pose as translation
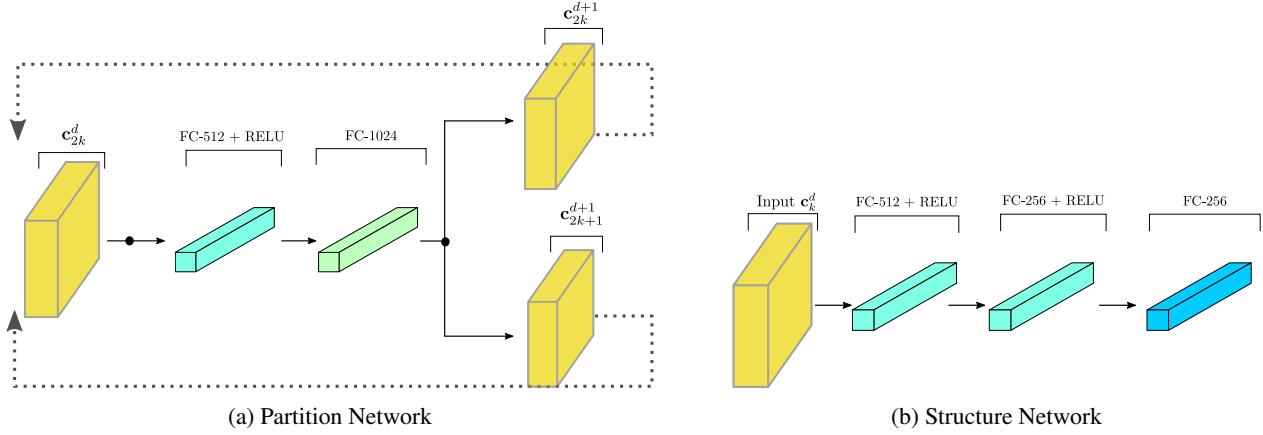
(a) Partition Network

(b) Structure Network

Figure 14: **Network Architecture Overview.** The *partition network* (14a) is simply one hidden layer fully connected network with RELU non linearity. The gray dotted lines indicate the recursive partition of the feature representation. Similarly, the *structure network* (14b) consists of two fully connected layers followed by RELU non linearity.



(a) Shape

(b) Size

(c) Translation

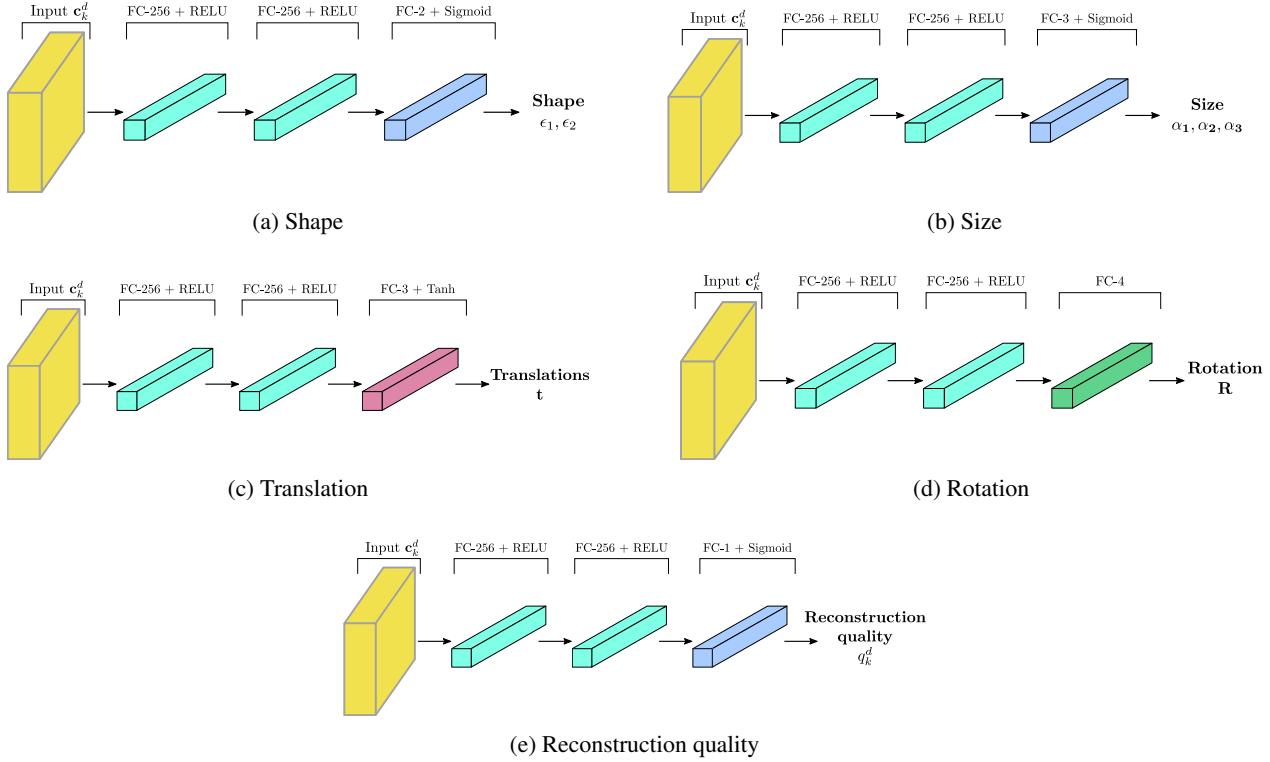(d) Rotation

(e) Reconstruction quality

Figure 15: **Geometry Network.** We detail the specifics of each regressor for predicting the primitive parameters $\lambda_k^d$ and the reconstruction quality $q_k^d$.

$\mathbf{t}$ and rotation $\mathbf{R}$) in addition to the reconstruction quality $q_k^d$. Fig. 15 presents the details of the implementation of each regressor.

## B.2. Training

In all our experiments, we use the Adam optimizer [28] with learning rate 0.0001 and no weight decay. For other hyper-parameters of Adam we use the PyTorch defaults. We train all models with a batch size of 32 for 40k iterations. We do not

perform any additional data augmentation. We weigh the loss terms of Eq. 9 in our main submission with 0.1, 0.01, 0.01 and 0.1 respectively, in order to enforce that during the first stages of training the network will focus primarily on learning the hierarchical decomposition of the 3D shape ($\mathcal{L}_s + \mathcal{L}_p$). In this way, after the part decomposition is learned, the network also focuses on the part geometries ($\mathcal{L}_r$). We also experimented with a two-stage optimization scheme, where we first learn the hierarchical part decomposition and then learn the hierarchical representation, but we observed that this made learning harder.

### B.3. Sampling Strategy

Sampling a point inside the target mesh has a probability proportional to the volume of the mesh. This yields bad reconstructions for thin parts of the object, such as legs of chairs and wings of aeroplanes. In addition, biasing the sampling towards the points inside the target mesh, results in worse reconstructions as also noted in [36]. To address the first issue (properly reconstructing thin parts), we use an unbalanced sampling distribution that, in expectation, results in sampling an equal number of points inside and outside the target mesh. To counter the second (biased sampling), we construct an unbiased estimator of the loss by weighing the per-point loss inversely proportionally to its sampling probability. We refer to our sampling strategy as *unbiased importance sampling*. Note that throughout all our experiments, we sample 10k points in the bounding box of the target mesh using our sampling strategy.

### B.4. Metrics

We evaluate our model and our baselines using the volumetric Intersection over Union (IoU) and the Chamfer-$L_1$ distance. Note that as our method does not predict a single mesh, we sample points from each primitive proportionally to its area, such that the total number of sampled points from all primitives is equal to 100k. For a fair comparison, we do the same for [53, 43]. Below, we discuss in detail the computation of the volumetric IoU and the Chamfer-$L_1$.

Volumetric IoU is defined as the quotient of the volume of the intersection of the target $S_{target}$ and the predicted $S_{pred}$ mesh and the volume of their union. We obtain unbiased estimates of the volume of the intersection and the union by randomly sampling 100k points from the bounding volume and determining if the points lie inside or outside the target / predicted mesh,

$$\text{IoU}(S_{pred}, S_{target}) = \frac{\mid V(S_{pred} \cap S_{target}) \mid}{\mid V(S_{pred} \cup S_{target}) \mid} \tag{18}$$

where $V(.)$ is a function that computes the volume of a mesh.

We obtain an unbiased estimator of the Chamfer-$L_1$ distance by sampling 100k points on the surface of the target $S_{target}$ and the predicted $S_{pred}$ mesh. We denote $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^{N}$ the set of points sampled on the surface of the target mesh and $\mathcal{Y} = \{\mathbf{y}_i\}_{i=1}^{N}$ the set of points sampled on the surface of the predicted mesh. We compute the Chamfer-$L_1$ as follows:

$$D_{\text{chamfer}}(\mathcal{X}, \mathcal{Y}) = \frac{1}{N} \sum_{\mathbf{x}_i \in \mathcal{X}} \min_{\mathbf{y}_j \in \cup \mathcal{Y}} \|\mathbf{x}_i - \mathbf{y}_j\| + \frac{1}{N} \sum_{\mathbf{y}_i \in \cup \mathcal{Y}} \min_{\mathbf{x}_j \in \mathcal{X}} \|\mathbf{y}_i - \mathbf{x}_j\| \tag{19}$$

The first term of (19) measures the *completeness* of the predicted shape, namely how far is on average the closest predicted point from a ground-truth point. The second term measures the *accuracy* of the predicted shape, namely how far on average is the closest ground-truth point from a predicted point.

To ensure a fair comparison with our baselines, we use the evaluation code of [36] for the estimation of both the Volumetric IoU and the Chamfer-$L_1$.

### B.5. Empirical Analysis of Loss Formulation

In this section, we investigate the impact of how various components of our model affect the performance on the single-image 3D reconstruction task.

#### B.5.1 Impact of Sampling Strategy

We first discuss how the sampling strategy affects the performance of our model. Towards this goal, we evaluate our model on the single-view 3D reconstruction task using three different sampling strategies: (a) uniform sampling in the bounding box that contains the target object (b) biased sampling (namely sampling an equal number of points inside and outside the target mesh without reweighing) and (c) unbiased importance sampling as described in Section B.3. All models are trained on the

|            | IoU       | Chamfer-$L_1$ |
|------------|-----------|---------------|
| Uniform    | 0.383     | 0.073         |
| Biased     | 0.351     | 0.041         |
| Importance | **0.491** | 0.073         |

(a) Influence of sampling strategy

|                | IoU       | Chamfer-$L_1$ |
|----------------|-----------|---------------|
| Importance 2k  | 0.370     | 0.074         |
| Importance 5k  | 0.380     | 0.076         |
| Importance 10k | **0.491** | **0.073**     |

(b) Influence of number of sampled points.

Table 3: **Sampling Strategy.** We evaluate the performance of our model while varying the sampling scheme and the number of the sampled points inside the bounding box of the target mesh. We report the volumetric IoU (higher is better) and the Chamfer distance (lower is better) on the test set of the "chair category".

"chair" object category of ShapeNet using the same network architecture, the same number of sampled points ($N$ =10k) and the same maximum number of primitives ($D = 16$). The quantitative results on the test set of the "chair" category are shown in Table 3. We observe that the proposed importance sampling strategy achieves the best results in terms of IoU.

Furthermore, we also examine the impact of the number of sampled points on the performance of our model. In particular, we train our model on the "chair" category while varying the number of sampled points inside the bounding box that contains the target mesh. As expected, increasing the number of sampled points results in an improvement in reconstruction quality. We empirically found that sampling 10k points results in the best compromise between training time and reconstruction performance.

### B.5.2 Impact of Proximity loss

In this section, we explain empirically the vanishing gradient problem that emerges from the use of the sigmoid in the occupancy function of (17). To this end, we train two variants of our model, one with and without the proximity loss of Eq. 15, in our main submission. For this experiment, we train both variants on D-FAUST for the single image 3D reconstruction task. Both models are trained for a maximum number of 32 primitives and $s = 10$ and for the same number of iterations.
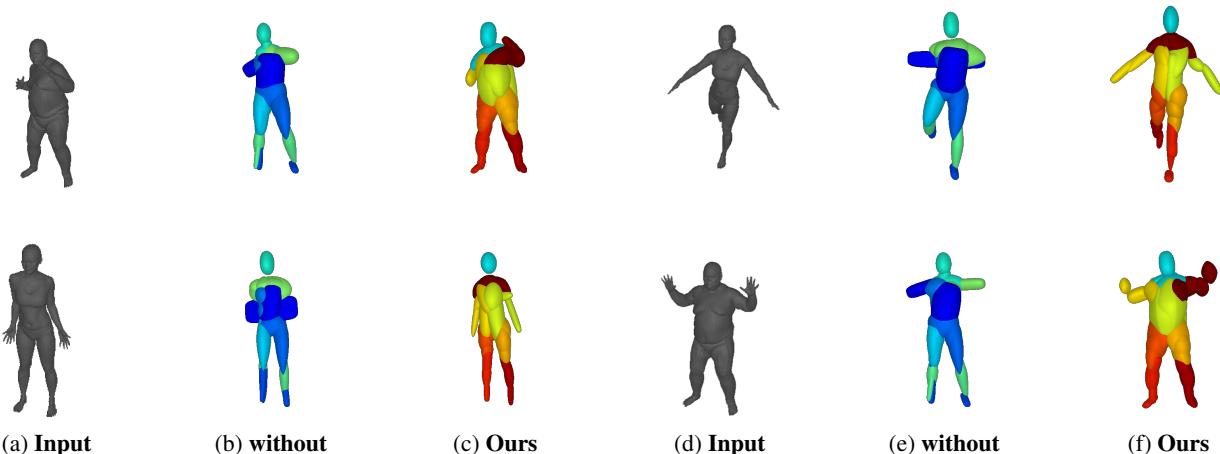


(a) **Input**  (b) **without**  (c) **Ours**  (d) **Input**  (e) **without**  (f) **Ours**

Figure 16: **Vanishing gradients.** We visualize the predictions of two two models, one trained with (**Ours**) and one **without** the proximity loss term. On the left, we visualize the input RGB image (a, d), in the middle the predictions without the proximity loss (b,c) and on the right the predictions of our model with this additional loss term.

|                         | IoU   | Chamfer-$L_1$ |
|-------------------------|-------|---------------|
| Ours w/o proximity loss | 0.605 | 0.171         |
| Ours                    | **0.699** | **0.098**  |

Table 4: **Proximity loss.** We investigate the impact of the proximity loss. We report the volumetric IoU and the Chamfer distance for two variants of our model, one with and without the proximity loss term.

Fig. 16 illustrates the predictions of both variants. We remark that the predictions of the model that was trained without the proximity loss are less accurate. Note that due to the vanishing gradient problem, the network is not able to properly "move" primitives and as a result, instead of reconstructing the hands of the humans using two or four primitives, the network uses only one. Interestingly, the reconstructions in some cases e.g. (e) do not even capture the human shape properly. However, even though the reconstruction quality is bad, the network is not able to fix it because the gradients of the reconstruction loss are small (even though the reconstruction loss itself is high). This is also validated quantitatively, as can be observed from Table 4.

## C. Additional Results on ShapeNet

In this section, we provide additional qualitative results on various object types from the ShapeNet dataset [6]. Furthermore, we also demonstrate the ability of our model to predict semantic hierarchies, where the same node is used for representing the same part of the object. We compare our model qualitatively with [43]. In particular, we train both models on the single-view 3D reconstruction task, using the same image renderings and train/test splits as [8]. Both methods are trained for a maximum number of 64 primitives. For our method, we empirically observed that a sharpness value $s = 10$ led to good reconstructions. Note that we do not compare qualitatively with [16, 9] as they do not provide code. Finally, we also compare our model with [53, 43] on the volumetric reconstruction task, where the input to all networks is a binary voxel grid. For a fair comparison, all models leverage the same feature encoder architecture proposed in [43].

In Fig. 18+19, we qualitatively compare our predictions with [43] for various ShapeNet objects. We observe that our model yields more accurate reconstructions compared to our baseline. Due to the use of the reconstruction quality $q_k^d$, our model dynamically decides whether a node should be split or not. For example, our model represents the phone in Fig. 18 (a) using one primitive (root node) and the phone in Fig. 18 (b), that consists of two parts, with two primitives. This can be also noted for the case of the displays Fig. 18 (g+j). For more complicated objects, such as aeroplanes, tables and chairs, our network uses more primitives to accurately capture the geometry of the target object. Note that for this experiment we set the threshold for $q_k^d$ to 0.8.

Our network associates the same node with the same part of the object, as it can be seen from the predicted hierarchies in Fig. 18+19. For example, for the displays the second primitive at the first depth level is used for representing the monitor of the display, for the aeroplanes the 4-th primitive in the second depth level is used for representing the front part of the aeroplanes.

## C.1. Volumetric Reconstruction

Our model is closely related to the works of Tulsiani et al. [53] and Paschalidou et al. [43]. Both [53, 43] were originally introduced using a binary occupancy grid as an input to their model, thus we also compare our model with [53, 43] using a voxelized input of size $32 \times 32 \times 32$. We evaluate the modelling accuracy of these three methods on the *animal* class of the ShapeNet dataset. To ensure a fair comparison, we use the feature encoder proposed in [53] for all three. A qualitative evaluation is provided in Fig. 17.

Our model yields more detailed reconstructions compared to [53, 43]. For example, in our reconstructions the legs of the animals are not connected and the tails better capture the geometry of the target shape. Again, we observe that our network predicts semantic hierarchies, where the same node is used for representing the same part of the animal.
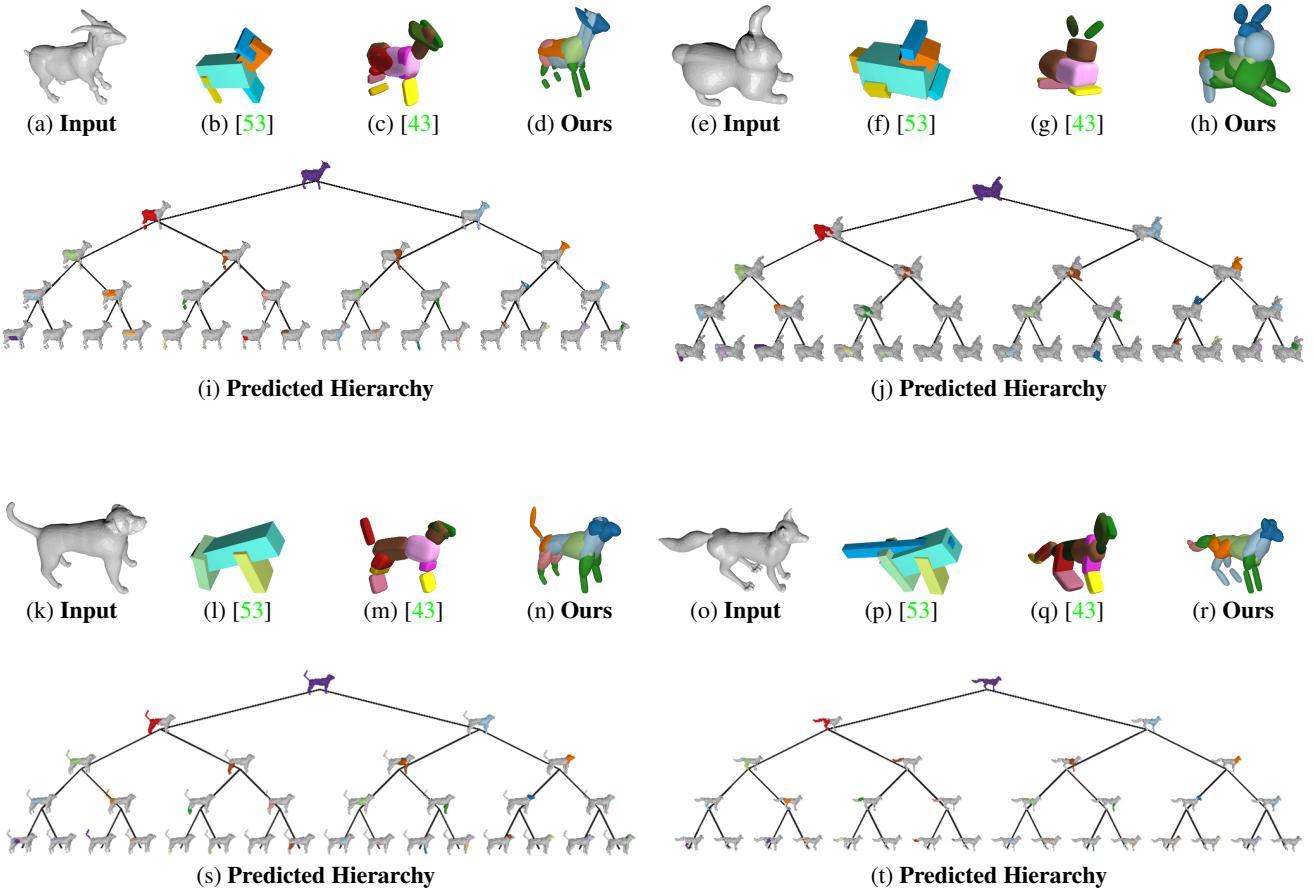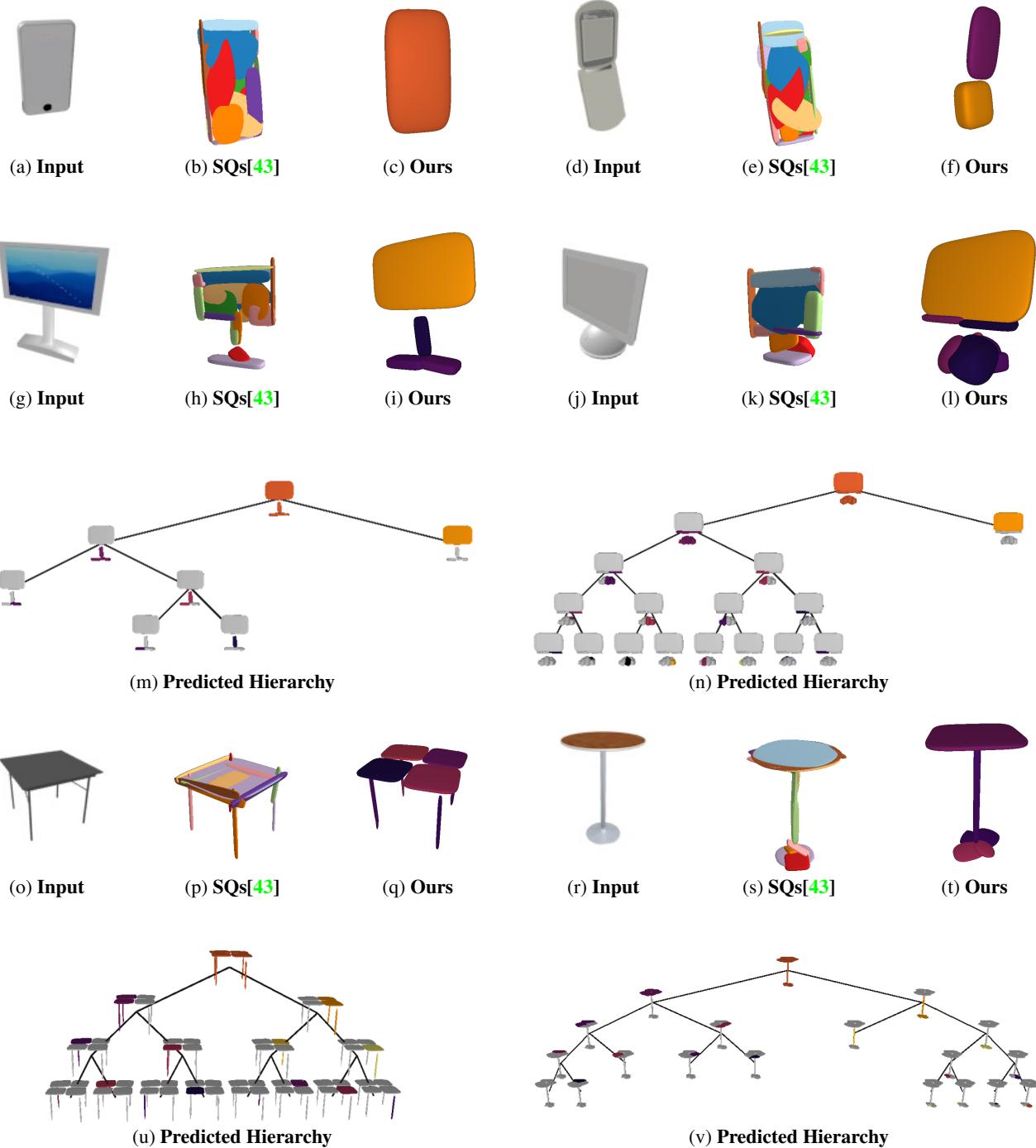
Figure 17: **Volumetric Reconstruction.** We note that our reconstructions are geometrically more accurate. In contrast to [43], our model yields reconstructions where the legs of the animals are not connected. Furthermore, our model accurately captures the ears and tails of the different animals.

## D. Additional Results on D-FAUST

In this section, we provide additional qualitative results on the D-FAUST dataset [4]. Furthermore, we also demonstrate that the learned hierarchies are indeed semantic as the same node is used to represent the same part of the human body. Similar to the experiment of Section 4.2 in our main submission, we evaluate our model on the single-view 3D reconstruction task, namely given a single *RGB image as an input*, our network predicts its geometry as a *tree of primitives as an output*. We compare our model with [43]. Both methods were trained for a maximum number of 32 primitives until convergence. For our method, we set the sharpness value $s = 10$.

In Fig. 20+22, we qualitatively compare our predictions with [43]. We remark that even though [43] is more parsimonious, our predictions are more accurate. For example, we note that our shape reconstructions capture the details of the muscles of the legs that are not captured in [43]. For completeness, we also visualize the predicted hierarchy up to the fourth depth level. Another interesting aspect of our model, which is also observed in [43, 53] is related to the semanticness of the learned hierarchies. We note that our model consistently uses the same node for representing the same part of the human body. For instance, node $(4, 15)$, namely the 15-th node at the 4-th depth level, consistently represents the right foot, whereas, node $(4, 12)$ represents the left foot. This is better illustrated in Fig. 21. In this figure, we only color the primitive associated with a particular node, for various humans, and we remark that the same primitive is used for representing the same body part.
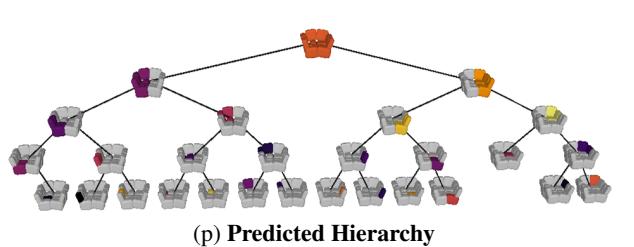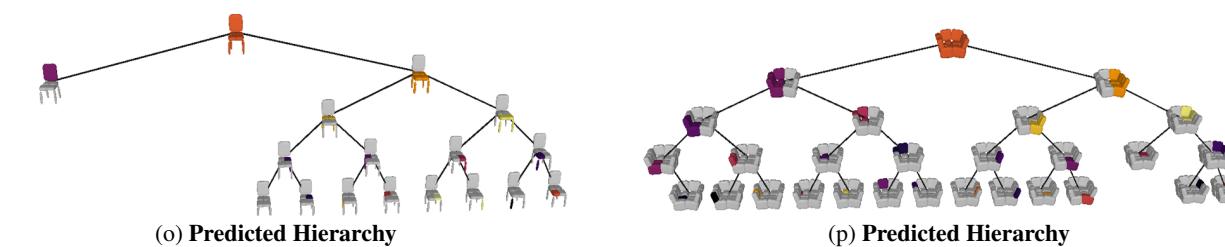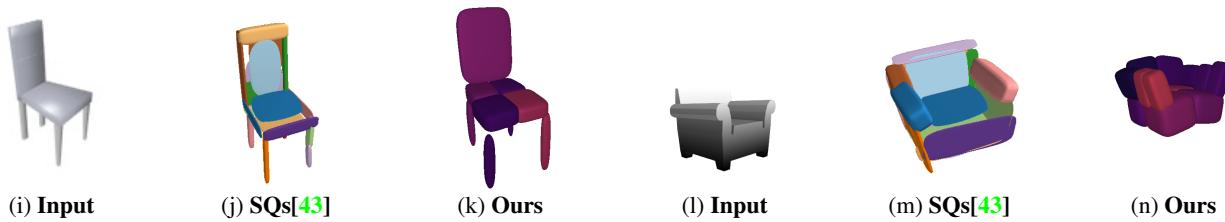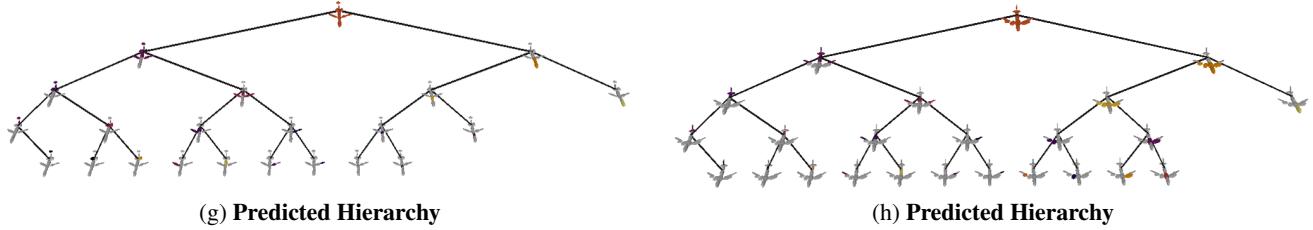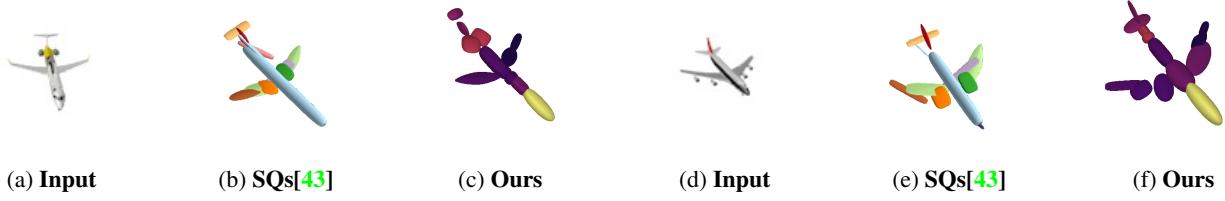
Finally, another interesting characteristic of our model is related to its ability to use less primitives for reconstructing humans, with smaller bodies. In particular, while the lower part of the human body is consistently represented with the same

(a) **Input**     (b) **SQs[43]**     (c) **Ours**     (d) **Input**     (e) **SQs[43]**     (f) **Ours**

(g) **Input**     (h) **SQs[43]**     (i) **Ours**     (j) **Input**     (k) **SQs[43]**     (l) **Ours**

(m) **Predicted Hierarchy**

(n) **Predicted Hierarchy**

(o) **Input**     (p) **SQs[43]**     (q) **Ours**     (r) **Input**     (s) **SQs[43]**     (t) **Ours**

(u) **Predicted Hierarchy**

(v) **Predicted Hierarchy**

Figure 18: **Single Image 3D Reconstruction on ShapeNet.** We visualize the predictions of our model on various ShapeNet objects and compare to [43]. For objects that are represented with more than two primitives, we also visualize the predicted hierarchy.
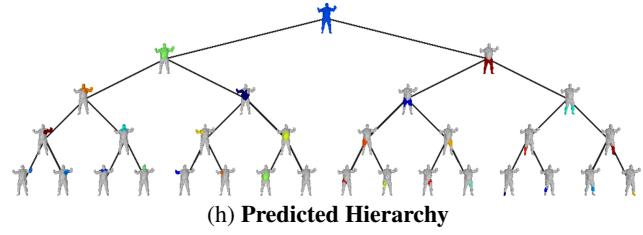
set of primitives, the upper part can be represented with less depending on the size and the articulation of the human body. This is illustrated in Fig. 22, where we visualize the predictions of our model for such scenarios.

Below, we provide the full hierarchies of the results on D-FAUST from our main submission.

(a) **Input**     (b) **SQs[43]**     (c) **Ours**     (d) **Input**     (e) **SQs[43]**     (f) **Ours**

(g) **Predicted Hierarchy**        (h) **Predicted Hierarchy**

(i) **Input**     (j) **SQs[43]**     (k) **Ours**     (l) **Input**     (m) **SQs[43]**     (n) **Ours**

(o) **Predicted Hierarchy**        (p) **Predicted Hierarchy**

Figure 19: **Single Image 3D Reconstruction on ShapeNet.** We visualize the predictions of our model on various ShapeNet objects and compare to [43]. For objects that are represented with more than two primitives, we also visualize the predicted hierarchy.
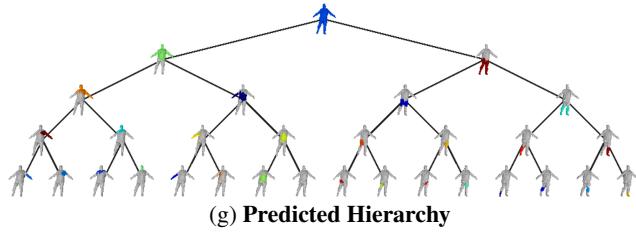
(a) **Input**   (b) **SQs[43]**   (c) **Ours**   (d) **Input**   (e) **SQs[43]**   (f) **Ours**

(g) **Predicted Hierarchy**   (h) **Predicted Hierarchy**

(i) **Input**   (j) **SQs[43]**   (k) **Ours**   (l) **Input**   (m) **SQs[43]**   (n) **Ours**

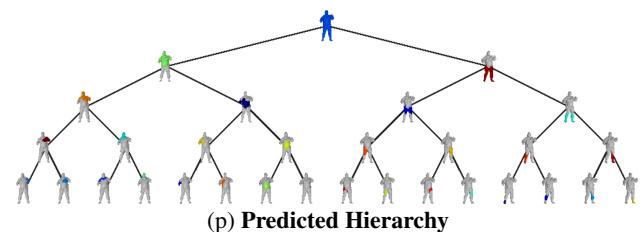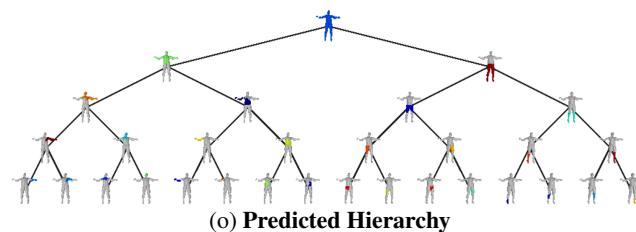(o) **Predicted Hierarchy**   (p) **Predicted Hierarchy**

Figure 20: **Qualitative Results on D-FAUST.** Our network learns semantic mappings of body parts across different body shapes and articulations while being geometrical more accurate compared to [43].

(a) Node $(4, 0)$

(b) Node $(3, 3)$

(c) Node $(4, 3)$
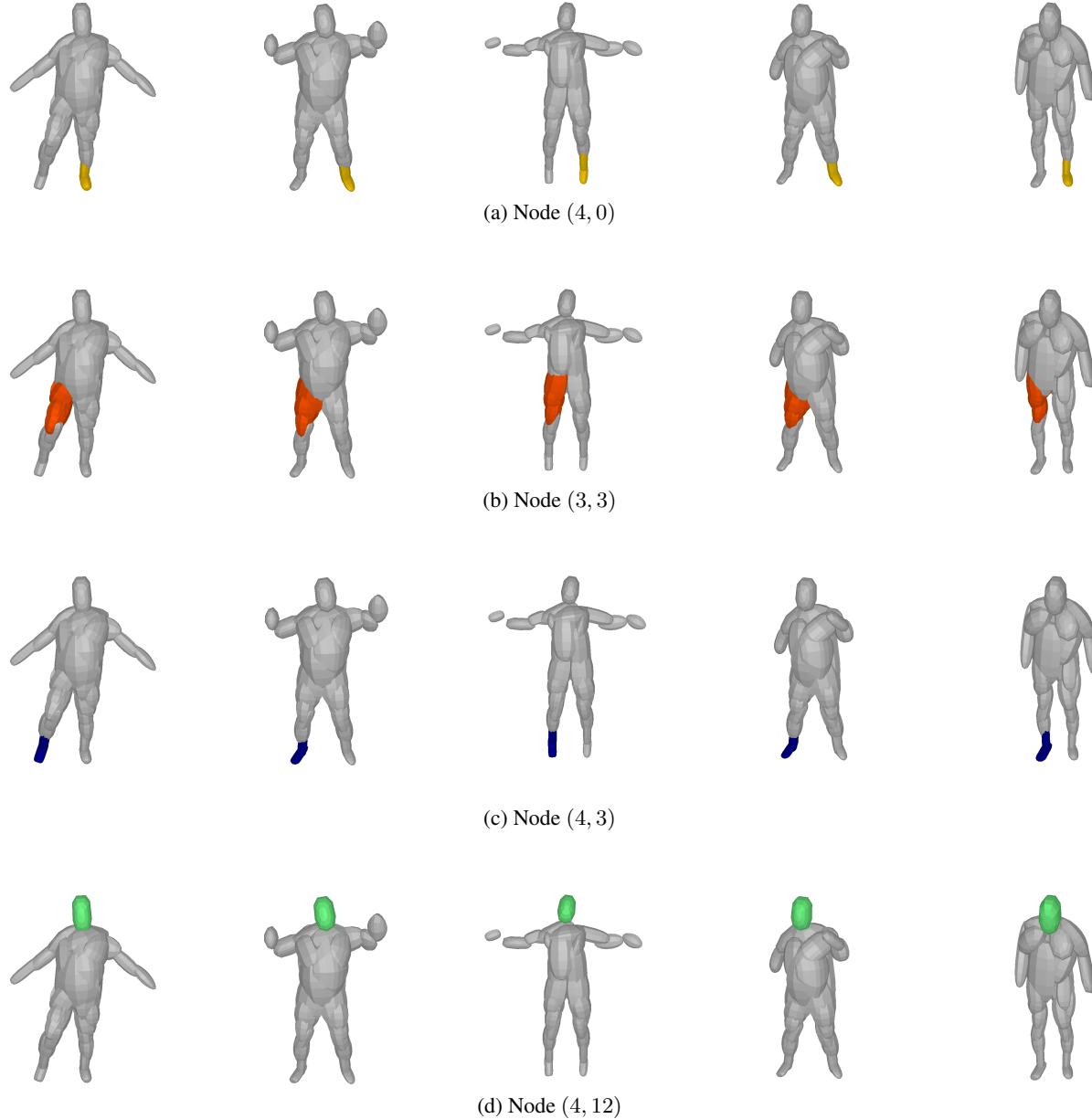
(d) Node $(4, 12)$

Figure 21: **Semantic Predictions on D-FAUST.** To illustrate that our model indeed learns semantic hierarchical layouts of parts, here we color a specific node of the tree for various humans and we observe that it consistently corresponds to the same body part.
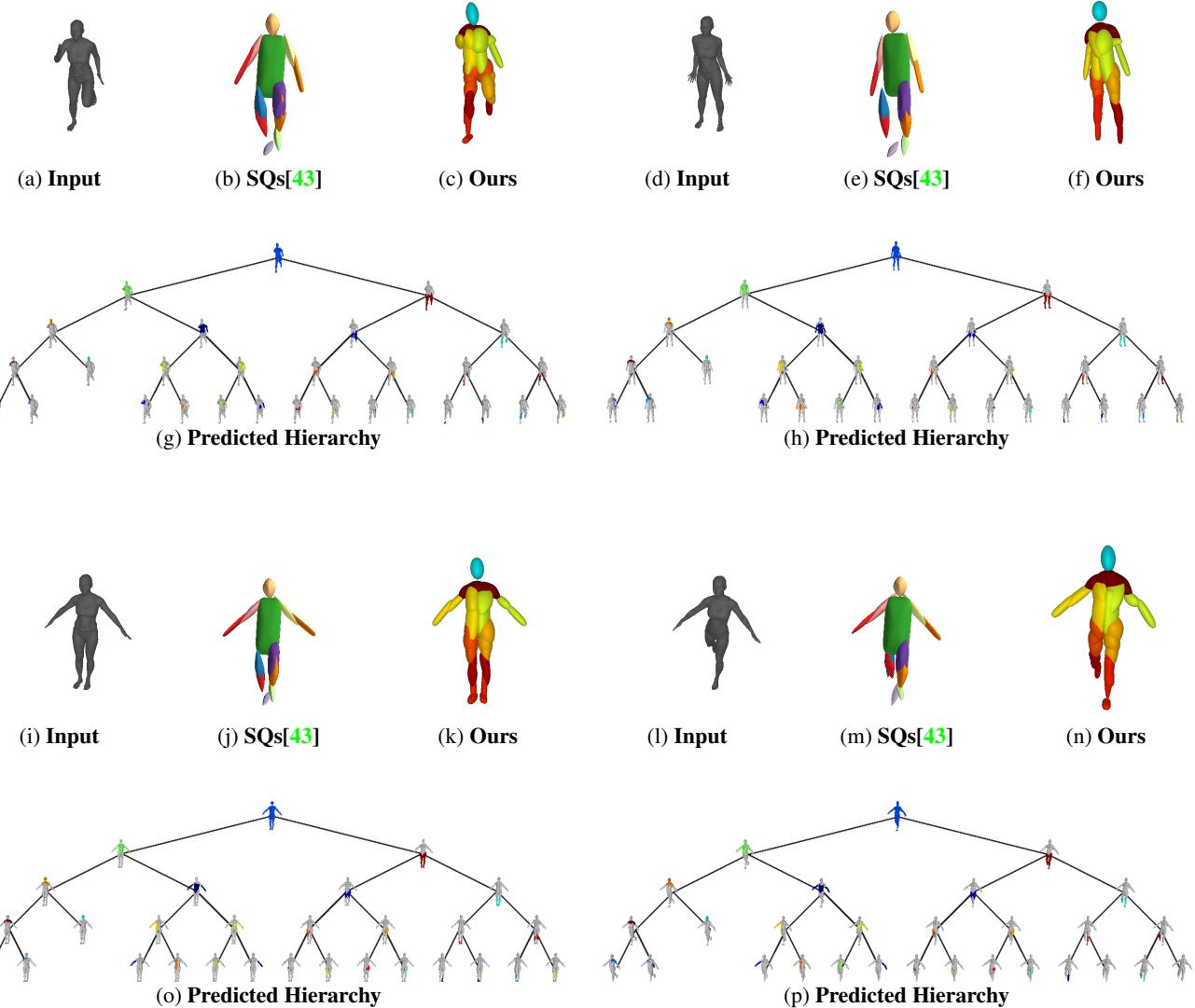
(a) **Input**      (b) **SQs[43]**      (c) **Ours**      (d) **Input**      (e) **SQs[43]**      (f) **Ours**

(g) **Predicted Hierarchy**             (h) **Predicted Hierarchy**

(i) **Input**      (j) **SQs[43]**      (k) **Ours**      (l) **Input**      (m) **SQs[43]**      (n) **Ours**

(o) **Predicted Hierarchy**             (p) **Predicted Hierarchy**

Figure 22: **Qualitative Results on D-FAUST.** Our network learns semantic mappings of body parts across different body shapes and articulations. Note that the network predicts less primitives for modelling the upper part of the human body.
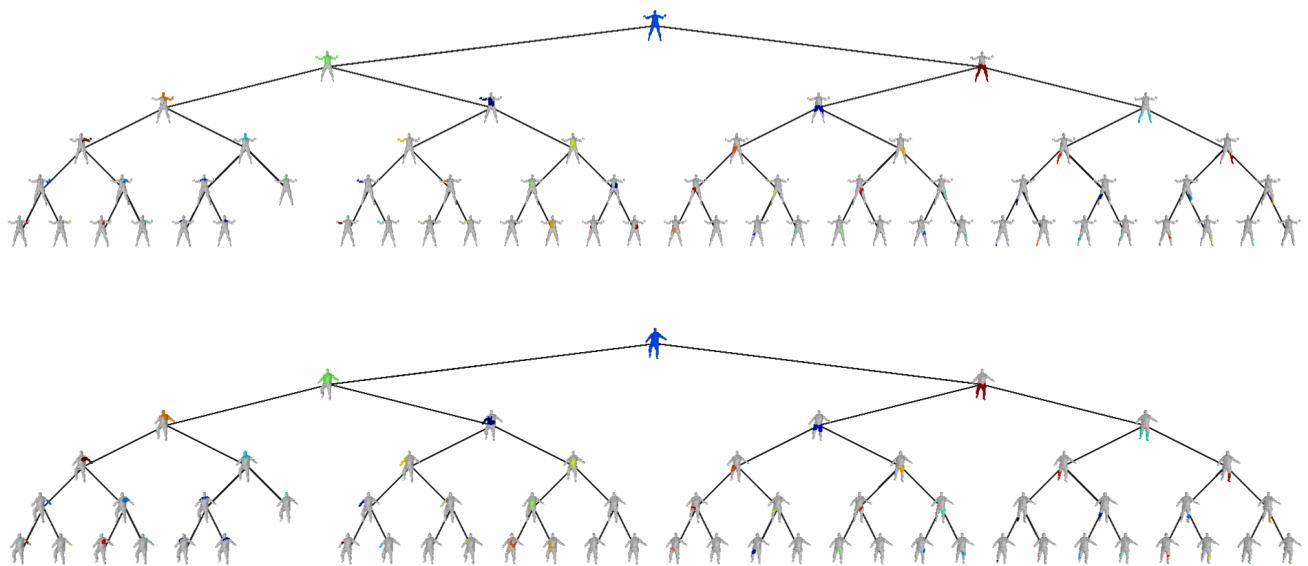
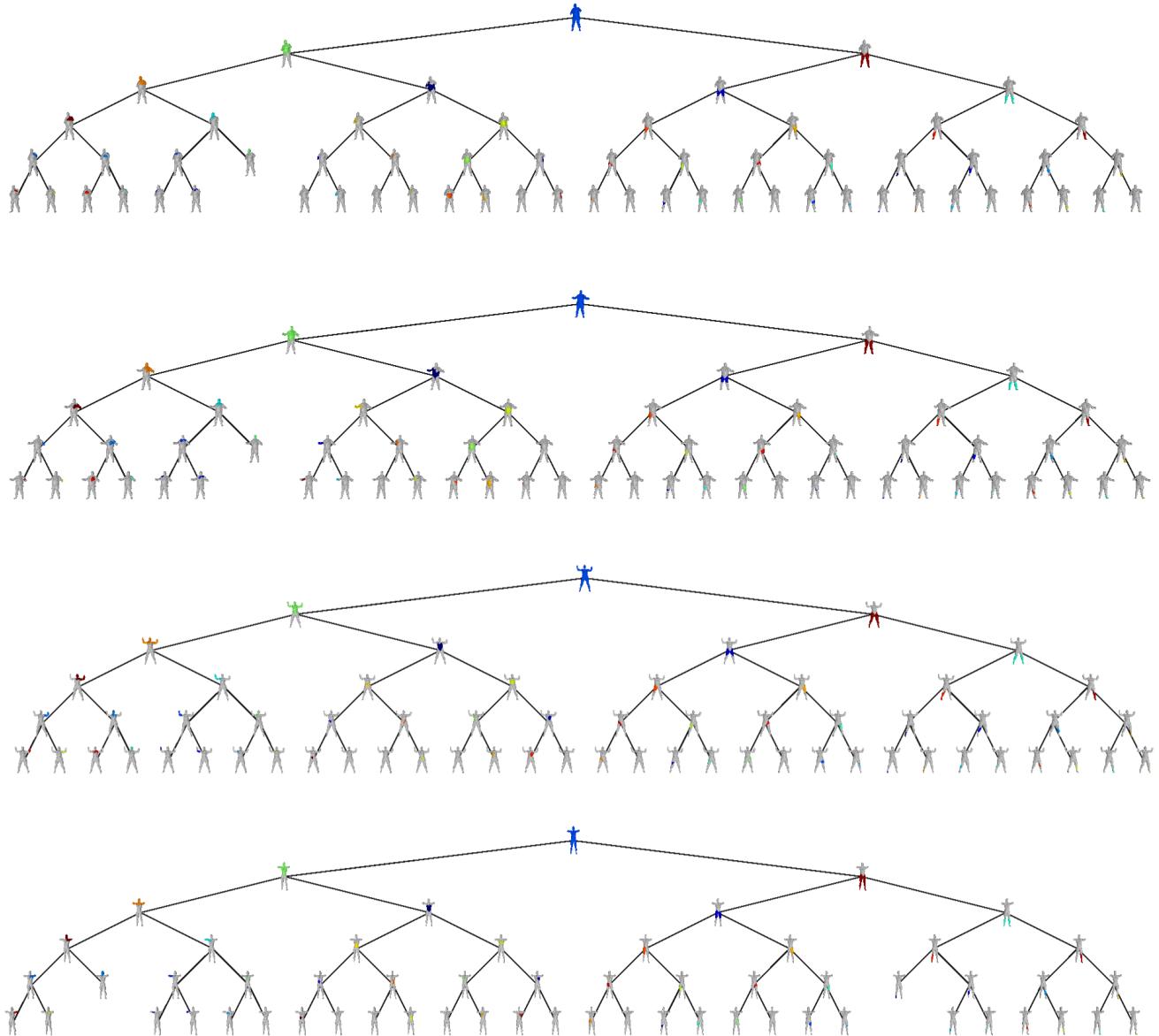Figure 23: **Full hierarchies of Figure 6 in our main submission.** Please zoom-in for details

Figure 24: **Full hierarchies of Figure 8 in our main submission.** Please zoom-in for details.