

Local Deep Implicit Functions for 3D Shape

Kyle Genova^{1,2} Forrester Cole² Avneesh Sud² Aaron Sarna² Thomas Funkhouser^{1,2}

¹Princeton University ²Google Research

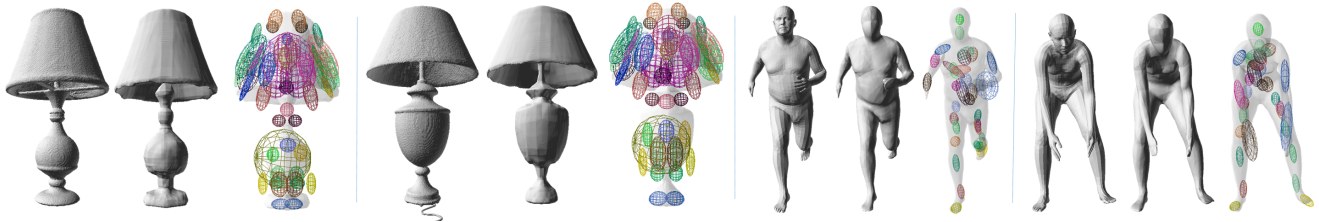


Figure 1. This paper introduces Local Deep Implicit Functions, a 3D shape representation that decomposes an input shape (mesh on left in every triplet) into a structured set of shape elements (colored ellipses on right) whose contributions to an implicit surface reconstruction (middle) are represented by latent vectors decoded by a deep network. Project video and website at ldif.cs.princeton.edu.

Abstract

The goal of this project is to learn a 3D shape representation that enables accurate surface reconstruction, compact storage, efficient computation, consistency for similar shapes, generalization across diverse shape categories, and inference from depth camera observations. Towards this end, we introduce Local Deep Implicit Functions (LDIF), a 3D shape representation that decomposes space into a structured set of learned implicit functions. We provide networks that infer the space decomposition and local deep implicit functions from a 3D mesh or posed depth image. During experiments, we find that it provides 10.3 points higher surface reconstruction accuracy (F-Score) than the state-of-the-art (OccNet), while requiring fewer than 1% of the network parameters. Experiments on posed depth image completion and generalization to unseen classes show 15.8 and 17.8 point improvements over the state-of-the-art, while producing a structured 3D representation for each input with consistency across diverse shape collections.

1. Introduction

Representing 3D shape is a fundamental problem with many applications, including surface reconstruction, analysis, compression, matching, interpolation, manipulation, and visualization. For most vision applications, a 3D representation should support: (a) reconstruction with accurate

surface details, (b) scalability to complex shapes, (c) support for arbitrary topologies, (d) generalizability to unseen shape classes, (e) independence from any particular application domain, (f) encoding of shape priors, (g) compact storage, and (h) computational efficiency.

No current representation has all of these desirable properties. Traditional explicit 3D representations (voxels, meshes, point clouds, etc.) provide properties (a-e) above. They can represent arbitrary shapes and any desired detail, but they don't encode shape priors helpful for efficient storage, 3D completion, and reconstruction tasks. In contrast, learned representations (latent vectors and deep network decoders) excel at representing shapes compactly with low-dimensional latent vectors and encoding shape priors in network weights, but they struggle to reconstruct details for complex shapes or generalize to novel shape classes.

Most recently, deep implicit functions (DIF) have been shown to be highly effective for reconstruction of individual objects [24, 7, 27, 45]. They represent an input observation as a latent vector \mathbf{z} and train a neural network to estimate the inside/outside or signed-distance function $f(\mathbf{x}, \mathbf{z})$ given a query location \mathbf{x} in 3D space. This approach achieves state of the art results for several 3D shape reconstruction tasks. However, they use a single, fixed-length latent feature vector to represent the entirety of all shapes and they evaluate a complex deep network to evaluate the implicit function for every position \mathbf{x} . As a result, they support limited shape complexity, generality, and computational efficiency.

Meanwhile, new methods are emerging for learning to infer structured decomposition of shapes [39, 13]. For example, [13] recently proposed a network to encode shapes into Structured Implicit Functions (SIF), which represents an implicit function as a mixture of local Gaussian functions. They showed that simple networks can be trained to decompose diverse collections of shapes consistently into SIFs, where the local shape elements inferred for one shape (e.g., the leg of a chair) correspond to semantically similar elements for others (e.g., the leg of a table). However, they did not use these structured decompositions for accurate shape reconstruction due to the limited shape expressivity of their local implicit functions (Gaussians).

The key idea of this paper is to develop a pipeline that can learn to infer *Local Deep Implicit Functions*, (“LDIF”, Figure 1). An LDIF is a set of local DIFs that are arranged and blended according to a SIF template. The representation is similar to SIF in that it decomposes a shape into a set of overlapping local regions represented by Gaussians. However, it also associates a latent vector with each local region that can be decoded with a DIF to produce finer geometric detail. Alternately, LDIF is similar to a DIF in that it encodes a shape as a latent vector that can be evaluated with a neural network to estimate the inside/outside function $f(\mathbf{x}, \mathbf{z})$ for any location \mathbf{x} . However, the LDIF latent vector is decomposed into parts associated with local regions of space (SIF Gaussians), which makes it more scalable, generalizable, and computationally efficient.

In this paper, we not only propose the LDIF representation, but we also provide a common system design that works effectively for 3D autoencoding, depth image completion, and partial surface completion. First, we propose to use DIF to predict local functions that are *residuals* with respect to the Gaussian functions predicted by SIF – this choice simplifies the task of the DIF, as it must predict only fine details rather than the overall shape within each shape element. Second, we propose to use the SIF decomposition of space to focus the DIF encoder on local regions by gathering input 3D points within each predicted shape element and encoding them with PointNet [30]. Finally, we investigate several significant improvements to SIF (rotational degrees of freedom, symmetry constraints, etc.) and simplifications to DIF (fewer layers, smaller latent codes, etc.) to improve the LDIF representation. Results of ablation studies show that each of these design choices provides significant performance improvements over alternatives. In all, LDIF achieves 10-15 points better F-Score performance on shape reconstruction benchmarks than the state-of-the-art [24], with fewer than 1% of the network parameters.

2. Related Work

Traditional Shape Representations: There are many existing approaches for representing shape. In computer

graphics, some of the foundational representations are meshes [2], point clouds [11], voxel grids [11], and implicit surfaces [31, 3, 4, 25, 26, 44]. These representations are popular for their simplicity and ability to operate efficiently with specialized hardware. However, they lack two important properties: they do not leverage a shape prior, and they can be inefficient in their expressiveness. Thus, traditional surface reconstruction pipelines based on them, such as Poisson Surface Reconstruction [19], require a substantial amount of memory and computation and are not good for completing unobserved regions.

Learned Shape Representations: To leverage shape priors, shape reconstruction methods began representing shape as a learned feature vector, with a trained decoder to a mesh [36, 14, 41, 15, 18], point cloud [10, 21, 46], voxel grid [9, 43, 5, 42], or octree [37, 33, 32]. Most recently, representing shape as a vector with an implicit surface function decoder has become popular, with methods such as Oc-cNet [24], ImNet [7], DeepSDF [27], and DISN [45]. These methods have substantially improved the state of the art in shape reconstruction and completion. However, they do not scale or generalize very well because the fundamental representation is a single fixed-length feature vector representing a shape globally.

Structured Shape Representations: To improve scalability and efficiency, researchers have introduced structured representations that encode the repeated and hierarchical nature of shapes. Traditional structured representations include scene graphs [11], CSG trees [11], and partition of unity implicits [26], all of which represent complex shapes as the composition of simpler ones. Learned representations include SDM-Net [12], GRASS [20], CSGNet [35], Volumetric Primitives [39], Superquadrics [28], and SIF [13]. These methods can decompose shapes into simpler ones, usually with high consistency across shapes in a collection. However, they have been used primarily for shape analysis (e.g. part decomposition, part-aware correspondence), not for accurate surface reconstruction or completion. Concurrent work [8, 17] adds a voxel grid structure to deep implicit functions. This helps preserve local detail but does not take advantage of consistent shape decomposition.

3. Local Deep Implicit Functions

In this paper, we propose a new 3D shape representation, Local Deep Implicit Functions (LDIF). The LDIF is a function that can be used to classify whether a query point \mathbf{x} is inside or outside a shape. It is represented by a set of N shape elements, each parameterized by 10 analytic shape variables θ_i and M latent shape variables \mathbf{z}_i :

$$\text{LDIF}(\mathbf{x}, \Theta, \mathbf{Z}) = \sum_{i \in [N]} g(\mathbf{x}, \theta_i)(1 + f(\mathbf{x}, \mathbf{z}_i)) \quad (1)$$

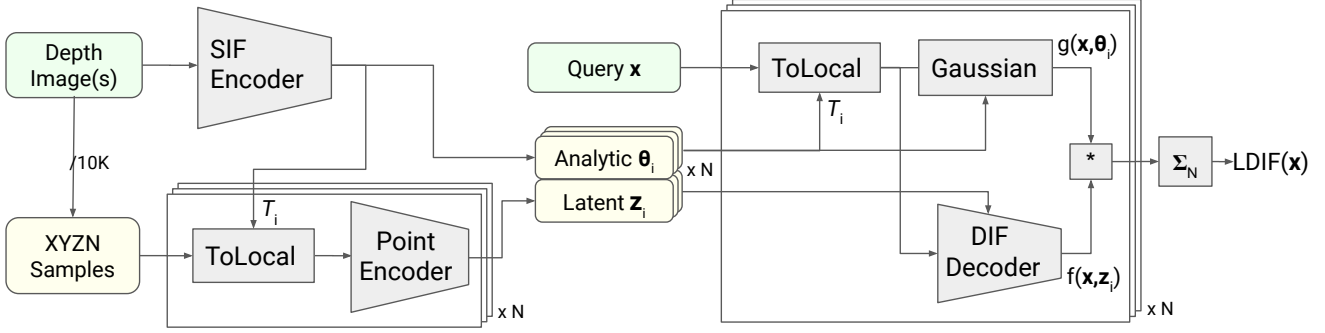


Figure 2. **Network architecture.** Our system takes in one or more posed depth images and outputs an LDIF function that can be used to classify inside/outside for any query point \mathbf{x} . It starts with a SIF encoder to extract a set of overlapping shape elements, each defined by a local Gaussian region of support parameterized by θ_i . It then extracts sample points/normals from the depth images and passes them through a PointNet encoder for each shape element to produce a latent vector \mathbf{z}_i . A local decoder network is used to decode each \mathbf{z}_i to produce an implicit function $f_i(\mathbf{x}, \mathbf{z}_i)$, which is combined with the local Gaussian function $g(\mathbf{x}, \theta_i)$ and summed with other shape elements to produce the output function $\text{LDIF}(\mathbf{x})$.

where $g(\mathbf{x}, \theta_i)$ is a local analytic implicit function and $f(\mathbf{x}, \mathbf{z}_i)$ is a deep implicit function. Intuitively, g provides a density function that defines a coarse shape and region of influence for each shape element, and f provides the shape details that cannot be represented by g .

Like a typical deep implicit function, our LDIF represents a 3D shape as an isocontour of an implicit function decoded with a deep network conditioned on predicted latent variables. However, LDIF replaces the (possibly long) single latent code of a typical DIF with the concatenation of N pairs of analytic parameters θ_i and short latent codes \mathbf{z}_i – i.e., the global implicit function is decomposed into the sum of N local implicit functions. This key difference helps it to be more accurate, efficient, consistent, scalable, and generalizable (see Section 6).

Analytic shape function. The analytic shape function g defines a coarse density function and region of influence for each shape element. Any simple analytic implicit function with local support would do. We use an oriented, anisotropic, 3D Gaussian:

$$g(\mathbf{x}, \theta_i) = c_i e^{-\frac{\|T_i \mathbf{x}\|^2}{2}} \quad (2)$$

where the parameter vector θ_i consists of ten variables: one for a scale constant c_i , three for a center point \mathbf{p}_i , three radii \mathbf{r}_i , and three Euler angles \mathbf{e}_i (this is the same parameterization as [13], except with 3 additional DoFs for rotation). The last 9 variables imply an affine transformation matrix T_i that takes a point \mathbf{x} from object space coordinates to the local isotropic, oriented, centered coordinate frame of the shape element.

Deep shape function. The deep implicit function f defines local shape details within a shape element by modulating g (one f function is shared by all shape elements). To compute f , we use a network architecture based on Occupancy

Networks [24]. As in the original OccNet, ours is organized as a fully-connected network conditioned on the latent code \mathbf{z}_i and trained using conditional batch normalization. However, one critical difference is that we transform the point \mathbf{x} by T_i before feeding it to the network. Another critical difference is that f_i only modulates the local implicit function g_i , rather than predicting an entire, global function. As a result, our local decoder has fewer network layers (9 vs. 33), shorter latent codes (32 vs. 256), and many fewer network parameters (8.6K vs 2M) than the original OccNet, and still achieves higher overall accuracy (see Section 6).

Symmetry constraints. For shape collections with man-made objects, we constrain a subset of the shape elements (half) to be symmetric with respect to a selected set of transformations (reflection across a right/left bisecting plane). These “symmetric” shape elements are evaluated twice for every point query, once for \mathbf{x} and once for $S\mathbf{x}$, where S is the symmetry transformation. In doing so, we effectively increase the number of shape elements without having to compute/store extra parameters for them. Adding partial symmetry encourages the shape decomposition to match global shape properties common in many shape collections and gives a boost to accuracy (Table 4).

4. Processing Pipeline

The processing pipeline for computing an LDIF is shown in Figure 2. All steps of the pipeline are differentiable and trained end-to-end. At inference time, the input to the system is a 3D surface or depth image, and the output is a set of shape element parameters Θ and latent codes \mathbf{Z} for each of N overlapping local regions, which can be decoded to predict inside/outside for any query location \mathbf{x} . Complete surfaces can be reconstructed for visualization by evaluating $\text{LDIF}(\mathbf{x})$ at points on a regular grid and running Marching

Cubes [23].

The exact configuration of the encoder architecture varies with input data type. We encode a **3D mesh** by first rendering a stack of 20 depth images at 137 x 137 resolution from a fixed set of equally spaced views surrounding the object. We then give the depth images to an early-fusion ResNet50 [16] to regress the shape element parameters Θ . Meanwhile, we generate a set of 10K points with normals covering the whole shape by estimating normals from the depth image(s) and unprojecting randomly selected pixels to points in object space using the known camera parameters. Then, for each shape element, we select a sampling of 1K points with normals within the region of influence defined by the predicted analytic shape function, and pass them to a PointNet [30] to generate the latent code \mathbf{z}_i . Alternatively, we could have encoded 3D input surfaces with CNNs based on mesh, point, or voxel convolutions, but found this processing pipeline to provide a good balance between detail, attention, efficiency, and memory. In particular, since the local geometry of every shape element is encoded independently with a PointNet, it is difficult for the network to “memorize” global shapes and it therefore generalizes better.

We encode a **depth image** with known camera parameters by first converting it into a 3 channel stack of 224 x 224 images representing the XYZ position of every pixel in object coordinates. We then feed those channels into a ResNet50 to regress the shape element parameters Θ , and we regress the latent codes \mathbf{Z} for each shape element using the same process as for 3D meshes.

4.1. Training Losses

The pipeline is trained with the following loss L :

$$L(\Theta, \mathbf{Z}) = w_P L_P(\Theta, \mathbf{Z}) + w_C L_C(\Theta) \quad (3)$$

Point Sample Loss L_P . The first loss L_P measures how accurately the LDIF(\mathbf{x}) predicts inside/outside of the ground-truth shape. To compute it, we sample 1024 points near the ground truth surface (set \mathcal{S}) and 1024 points uniformly at random in the bounding box of the shape (set \mathcal{U}). We combine them with weights $w_i \in \{w_S, w_U\}$ to form set $C = \mathcal{U} \cup \mathcal{S}$. The near-surface points are computed using the sampling algorithm of [13]. We scale by a hyperparameter α , apply a sigmoid to the decoded value LDIF(\mathbf{x}), and then compute an L_2 loss to the ground truth indicator function $I(\mathbf{x})$ (see [13] for details):

$$L_P(\Theta, \mathbf{Z}) = \frac{1}{|C|} \sum_{\mathbf{x}_i \in C} w_i \|\text{sig}(\alpha \text{LDIF}(\mathbf{x}_i, \Theta, \mathbf{Z})) - I(\mathbf{x}_i)\|$$

Shape Element Center Loss L_C . The second loss L_C encourages the center of every shape element to reside within

the target shape. To compute it, we estimate a signed distance function on a low-res 32x32x32 grid \mathbf{G} for each training shape. The following loss is applied based on the grid value $G(\mathbf{p}_i)$ at the center \mathbf{p}_i of each shape element:

$$L_C(\Theta) = \begin{cases} \sum_{\theta_i \in \Theta} G(\mathbf{p}_i)^2 & G(\mathbf{p}_i) > \beta \\ 0 & G(\mathbf{p}_i) \leq \beta \end{cases}$$

Here, β is a threshold chosen to account for the fact that \mathbf{G} is coarse. It is set to half the width of a voxel cell in \mathbf{G} . This setting makes it a conservative loss: it says that when \mathbf{p}_i is definitely outside the ground truth shape, \mathbf{p}_i should be moved inside. L_C never penalizes a center that is within the ground truth shape.

It is also possible for the predicted center to lie outside the bounding box of \mathbf{G} . In this case, there is no gradient for L_C , so we instead apply the inside-bounding-box loss from [13] using the object-space bounds of \mathbf{G} .

5. Experimental Setup

We execute a series of experiments to evaluate the proposed LDIF shape representation, compare it to alternatives, study the effects of its novel components, and test it in applications. Except where otherwise noted, we use $N = 32$ shape elements and $M = 32$ dimensional latent vectors during all experiments.

Datasets. When not otherwise specified, experiments are run on the ShapeNet dataset [6]. We use the train and test splits from 3D-R²N² [9]. We additionally subdivide the train split to create an 85%, 5%, 10% train, validation, and test split. We pre-process the shapes to make them watertight using the depth fusion pipeline from Occupancy Networks [24]. We train models multi-class (all 13 classes together) and show examples only from the test split.

Metrics. We evaluate shape reconstruction results with mean intersection-over-union (IoU) [24], mean Chamfer distance [24], and mean F-Score [38] at $\tau = 0.01$. As suggested in [38], we find that IoU is difficult to interpret for low values, and Chamfer distance is outlier sensitive, and so we focus our discussions mainly on F-Scores.

Baselines. We compare most of our results to the two most related prior works: Occupancy Networks [24] (OccNet), the state-of-the-art in deep implicit functions, and Structured Implicit Functions [13] (SIF), the state-of-the-art in structural decomposition. We also compare to the AtlasNet autoencoder [15], which predicts meshes explicitly.

6. Experimental Evaluation

In this section, we report results of experiments that compare LDIF and baselines with respect to how well they satisfy desirable properties of a 3D shape representation.

Category	IoU (\uparrow)			Chamfer (\downarrow)				F-Score (\uparrow , %)			
	Occ.	SIF	Ours	Occ.	SIF	Atl.	Ours	Occ.	SIF	Atl.	Ours
airplane	77.0	66.2	91.2	0.16	0.44	0.17	0.10	87.8	71.4	85.1	96.9
bench	71.3	53.3	85.6	0.24	0.82	0.31	0.17	87.5	58.4	76.8	94.8
cabinet	86.2	78.3	93.2	0.41	1.10	0.81	0.33	86.0	59.3	71.5	92.0
car	83.9	77.2	90.2	0.61	1.08	0.70	0.28	77.5	56.6	74.2	87.2
chair	73.9	57.2	87.5	0.44	1.54	1.05	0.34	77.2	42.4	60.7	90.9
display	81.8	69.3	94.2	0.34	0.97	0.54	0.28	82.1	56.3	71.4	94.8
lamp	56.5	41.7	77.9	1.67	3.42	1.57	1.80	62.7	35.0	51.1	83.5
rifle	69.5	60.4	89.9	0.19	0.42	0.16	0.09	86.2	70.0	85.6	97.3
sofa	87.2	76.0	94.1	0.30	0.80	0.50	0.35	85.9	55.2	70.0	92.8
speaker	82.4	74.2	90.3	1.01	1.99	1.31	0.68	74.7	47.4	60.7	84.3
table	75.6	57.2	88.2	0.44	1.57	1.07	0.56	84.9	55.7	67.5	92.4
telephone	90.9	83.1	97.6	0.13	0.39	0.16	0.08	94.8	81.8	89.6	98.1
watercraft	74.7	64.3	90.1	0.41	0.78	0.41	0.20	77.3	54.2	74.4	93.2
mean	77.8	66.0	90.0	0.49	1.18	0.67	0.40	81.9	59.0	72.2	92.2

Table 1. **Autoencoder results.** Comparison of 3D-R²N² test set reconstruction errors for OccNet (“Occ.”) [24], SIF (“SIF”) [13], AtlasNet (“Atl.”) [15], and LDIF (“Ours”) autoencoders.

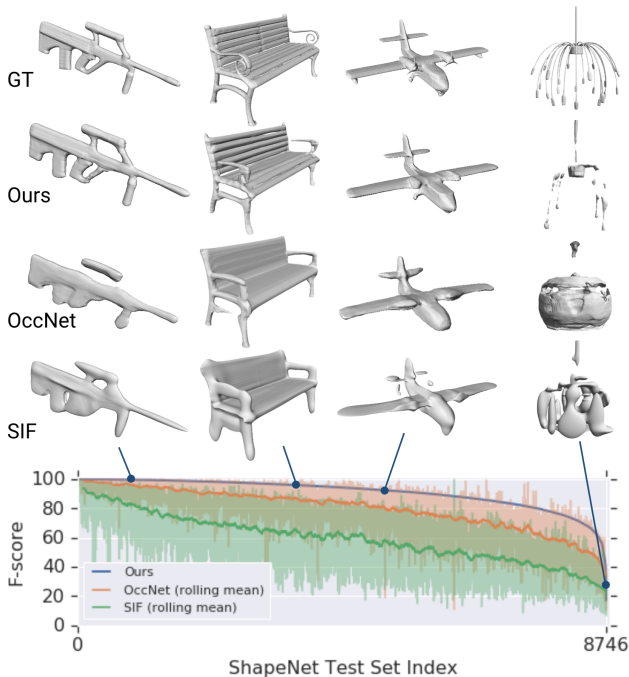


Figure 3. **Autoencoder examples.** F-scores for the test set (8746 shapes) are shown ordered by the LDIF F-score, with examples marked with their position on the curve. Our reconstructions (blue curve) are most accurate for 93% of shapes (exact scores shown faded). The scores of OccNet and SIF follow roughly the same curve as LDIF (rolling means shown bold), indicating shapes are similarly difficult for all methods. Solid shapes such as the rifle are relatively easy to represent, while shapes with irregular, thin structures such as the lamp are more difficult.

Accuracy. Our first experiment compares 3D shape representations in terms of how accurately they can encode/decode shapes. For each representation, we compare a 3D→3D autoencoder trained on the multiclass training data, use it to reconstruct shapes in the test set, and then

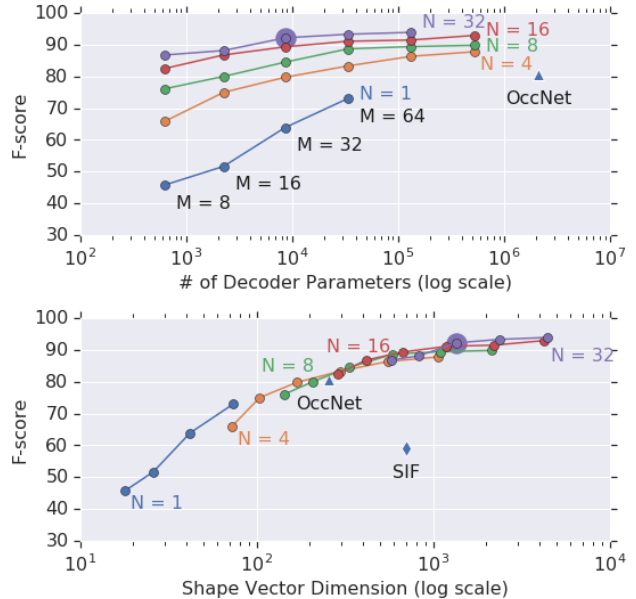


Figure 4. **Representation efficiency.** F-score vs. model complexity. Curves show varying M for constant N . Other methods marked as points. **Top:** F-score vs. count of decoder parameters. The $N = 32, M = 32$ configuration (large dot) reaches $>90\%$ F-score with $<1\%$ of the parameters of OccNet, and is used as the benchmark configuration in this paper. **Bottom:** F-score vs. shape vector dimension ($|\Theta| + |\mathbf{Z}|$ for DSIF). DSIF achieves similar reconstruction accuracy to OccNet at the same dimensionality, and can use additional dimensions to further improve accuracy.

evaluate how well the reconstructions match the originals (Table 1). LDIF’s mean F-Score is 92.2, 10.3 points higher than OccNet, 20.0 points higher than AtlasNet, and 33.2 points higher than SIF. A more detailed breakdown of the results appears in Figure 3, which shows the F-scores for all models in the test set – LDIF improves on OccNet’s score for 93% of test shapes. The increase in accuracy translates into a large qualitative improvement in results (shown above in Figure 3). For example, LDIF often reproduces better geometric details (e.g., back of the bench) and handles unusual part placements more robustly (e.g., handles on the rifle).

Efficiency. Our second experiment compares the efficiency of 3D shape representations in terms of accuracy vs. storage/computation costs. Since LDIF can be trained with different numbers of shape elements (N) and latent feature sizes (M), a family of LDIF representations is possible, each with a different trade-off between storage/computation and accuracy. Figure 4 investigates these trade-offs for several combinations of N and M and compares the accuracy of their autoencoders to baselines. Looking at the plot on the top, we see that LDIF provides more accurate reconstructions than baselines at every decoder size – our decoder with $N = 32$ and $M = 32$ is $0.004\times$ the size of OccNet and provides $1.13\times$ better F-Score. On the bottom, we see that LDIF performs comparably to OccNet and outperforms SIF

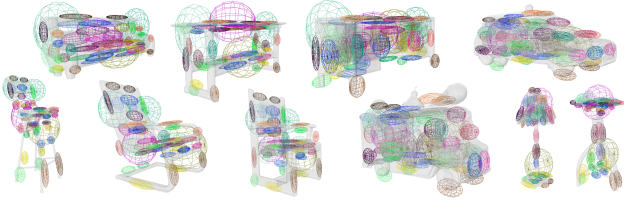


Figure 5. **Representation consistency.** Example shape decompositions produced by our model trained multi-class on 3D-R²N². Shape elements are depicted by their support ellipsoids and colored consistently by index. Note that the shape element shown in brown is used to represent the right-front leg of the chairs, tables, desks, and sofas, as well as the front-right wheel of the cars.

at the same number of bytes, despite having both deep and analytic parameters, and that it scales to larger embeddings.

Consistency. Our third experiment investigates the ability of LDIF to decompose shapes consistently into shape elements. This property was explored at length in [13] and shown to be useful for structure-aware correspondences, interpolations, and segmentations. While not the focus of this paper, we find qualitatively that the consistency of the LDIF representation is slightly superior to SIF, because the shape element symmetries and rotations introduced in this paper provide the DoFs needed to decompose shapes with fewer elements. On the other hand, the local DIFs are able to compensate for imperfect decompositions during reconstruction, which puts less pressure on consistency. Figure 5 shows qualitative results of the decompositions computed for LDIF. Please note the consistency of the colors (indicating the index of the shape element) across a broad range of shapes.

Generalizability. Our third experiment studies how well trained autoencoders generalize to handle unseen shape classes. To test this, we used the auto-encoders trained on 3D-R²N² classes and tested them without fine-tuning on a random sampling of meshes from 10 ShapeNet classes that were not seen during training. Table 2 shows that the mean F-Score for LDIF on these novel classes is 84.4, which is 17.8 points higher than OccNet and 41.4 points higher than SIF. Looking at the F-Score for every example in the bottom of Figure 6, we see that LDIF is better on 91% of examples. We conjecture this is because LDIF learns to produce consistent decompositions for a broad range of input shapes when trained multiclass, and because the local encoder network learns to predict shape details only for local regions. This two-level factoring of structure and detail seems to help LDIF generalize.

Domain-independence. Our fifth experiment investigates whether LDIF can be used in application domains beyond the man-made shapes found in ShapeNet. As one example, we trained LDIF without any changes to autoencode meshes of human bodies in a wide variety of poses sam-

Category	Chamfer (↓)			F-Score (↑, %)		
	SIF	OccNet	Ours	SIF	OccNet	Ours
bed	2.24	1.30	0.68	32.0	59.3	81.4
birdhouse	1.92	1.25	0.75	33.8	54.2	76.2
bookshelf	1.21	0.83	0.36	43.5	66.5	86.1
camera	1.91	1.17	0.83	37.4	57.3	77.7
file	0.71	0.41	0.29	65.8	86.0	93.0
mailbox	1.46	0.60	0.40	38.1	67.8	87.6
piano	1.81	1.07	0.78	39.8	61.4	82.2
printer	1.44	0.85	0.43	40.1	66.2	84.6
stove	1.04	0.49	0.30	52.9	77.3	89.2
tower	1.05	0.50	0.47	45.9	70.2	85.7
mean	1.48	0.85	0.53	43.0	66.6	84.4

Table 2. **Generalization to unseen classes.** Comparison of 3D reconstruction accuracy when 3D autoencoders are tested directly on ShapeNet classes not seen during training. Note that our method (LDIF) has a higher F-Score by 17.8 points.

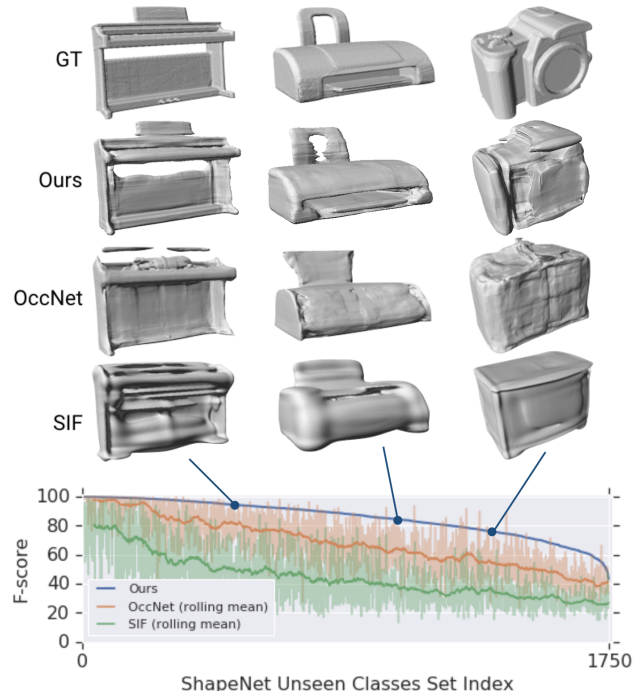


Figure 6. **Generalization examples.** Example shape reconstructions for piano, printer, and camera classes, which did not appear in the training data. F-score is plotted below ordered by LDIF score, similar to Figure 3. Our method (blue curve) achieves the best accuracy on 91% of the novel shapes.

pled from [40]. Specifically, we generated 5M meshes by randomly sampling SMPL parameters (CAESAR fits for shape, mocap sequence fits for pose). We used an 80%, 5%, 15% split for the train, val, and test sets, similar to [40], and measured the error of the learned autoencoder on the held-out test set. The challenge for this dataset is quite dif-



Figure 7. **Human body modeling.** Surface reconstructions and decompositions for 4 random SMPL [22] human meshes from the SURREAL [40] dataset. For each triple, from left to right: SMPL mesh, our reconstruction, our shape decomposition. These results demonstrate unsupervised correspondence between people in different poses as well as accurate reconstructions of organic shapes.

ferent than for ShapeNet – the autoencoder must be able to represent large-scale, non-rigid deformations in addition to shape variations. Our reconstructions achieve 93% mIOU compared to 85% mIOU for SIF. The results of LDIF reconstructions and the underlying SIF templates are shown in Figure 7. Despite a lack of supervision on pose or subject alignment, our approach reconstructs a surface close to the original and establishes coarse correspondences.

7. Applications

In this section, we investigate how the proposed LDIF representation can be used in applications. Although SIF (and similarly LDIF) has previously been shown useful for 3D shape analysis applications like structure-aware shape interpolation, surface correspondence, and image segmentation [13], we focus our study here on 3D surface reconstruction from partial observations.

7.1. 3D Completion from a Single Depth Image

Task. Reconstructing a complete 3D surface from a single depth image is an important vision task with applications in AR, robotics, etc. To investigate how LDIF performs on this task, we modified our network to take a single depth image as input (rather than a stack of 20) and trained it from scratch on depth images generated synthetically from random views of the 3D-R²N² split of shapes. The depth images were 512 x 512 to approximate the resolution of real depth sensors (though all CNN inputs are 224 x 224 due to memory restrictions). The depth images were rendered from view points sampled from all view directions and at variable distances to mimic the variety of scan poses. Each depth image was then converted to a three channel XYZ

Category	IoU (\uparrow)		Chamfer (\downarrow)		F-Score (\uparrow , %)	
	OccNet*	Ours	OccNet*	Ours	OccNet*	Ours
airplane	-	80.2	0.47	0.17	70.1	89.2
bench	-	70.9	0.70	0.39	64.9	81.9
cabinet	-	82.8	1.13	0.77	70.1	77.9
car	-	81.4	0.99	0.51	61.6	72.4
chair	-	70.6	2.34	1.02	50.2	69.6
display	-	82.4	0.95	0.62	62.8	80.0
lamp	-	62.1	9.91	2.15	44.1	66.4
rifle	-	81.5	0.49	0.14	66.4	92.3
sofa	-	81.4	1.08	0.83	61.2	71.7
speaker	-	80.2	3.50	1.48	52.4	67.3
table	-	73.5	2.49	1.14	66.7	78.0
telephone	-	92.3	0.35	0.19	86.1	92.0
watercraft	-	76.0	1.15	0.50	54.5	77.5
mean	-	78.1	1.97	0.76	62.4	78.2

Table 3. **Depth completion accuracy.** Our method (LDIF) provides better 3D surface completions than an OccNet* trained on our XYZ image inputs.

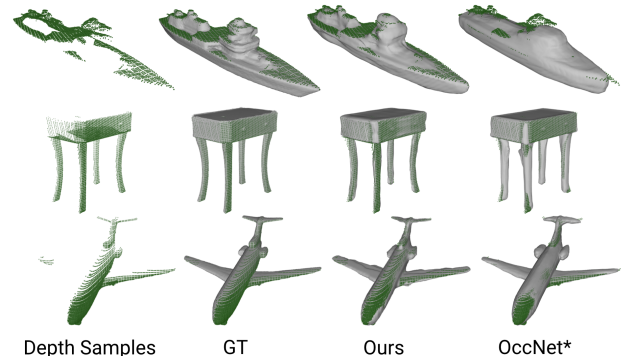


Figure 8. **Depth completion examples.** Visualizations of surfaces predicted from posed depth images (depicted by green points). Our method provides better details in both the observed and unobserved parts of the shape.

image using the known camera parameters.

Baseline. For comparison, we trained an OccNet network from scratch on the same data. Because the OccNet takes a point cloud rather than depth images, we train an XYZ image encoder network to regress the 256-D OccNet embedding. This OccNet* model provides an apples-to-apples baseline that isolates differences due only to the representation decoding part of the pipeline.

Results. Table 3 shows results of this 3D depth completion experiment. We find that the F-Score of LDIF is 15.8 points higher than OccNet* (78.2 vs. 62.4). Figure 8 highlights the difference in the methods qualitatively. As in the 3D case, we observe that LDIF’s local part encoders result in substantially better performance on hard examples.

Ablation study. To further understand the behavior of LDIF during depth completion, we ablate three components of our pipeline (Table 4). First, we verify that having local

Method	IoU (\uparrow)	Chamfer (\downarrow)	F-Score (\uparrow , %)
Full (D)	77.2	0.78	77.6
No PointNet	69.1	0.98	66.2
No Transform	71.9	1.80	71.9
No Symmetry	76.7	0.76	76.6

Table 4. **Depth completion ablation study.** Local PointNet encoders, camera transformations, and partial symmetry all improve performance. Independently and locally encoding the z_i with PointNet is particularly good for generalization (see Section 6).

pointnets to encode the local feature vectors is useful, rather than simply predicting them directly from the input image. Second, we show that providing an XYZ image as input to the network is much more robust than providing a depth image. Finally, we show that taking advantage of the explicit structure via partial symmetry improves results qualitatively and achieves the same quality with fewer degrees of freedom. The biggest of these differences is due to the PointNet encoding of local shape elements, which reduces the F-Score by 11.4 points if it is disabled.

7.2. Reconstruction of Partial Human Body Scans

Task. Acquisition of complete 3D surface scans for a diverse collection of human body shapes has numerous applications [1]. Unfortunately, many real world body scans have holes (Figure 9a), due to noise and occlusions in the scanning process. We address the task of learning to complete and beautify the partial 3D surfaces without any supervision or even a domain-specific template.

Dataset and baselines. The dataset for this experiment is CAESAR [34]. We use our proposed 3D autoencoder to learn to reconstruct an LDIF for every scan in the CAESAR dataset, and then we extract watertight surface from the LDIFs (using the splits from [29]). For comparisons, we do the same for SIF (another unsupervised method) and a non-rigid deformation fit of the S-SCAPE template [29].

Results. Figure 9 shows representative results. Note that LDIF captures high-frequency details missing in SIF reconstructions. Although the approach based on S-SCAPE provides better results, it requires a template designed specifically for human bodies as well as manual supervision (landmarks and bootstrapping), whereas LDIF is domain-independent and unsupervised. These results suggest that LDIF could be used for 3D reconstruction of other scan datasets where templates are not available.

8. Conclusion

Summary of research contributions: In this paper, we propose Local Deep Implicit Functions (LDIF), a new 3D representation that describes a shape implicitly as the sum

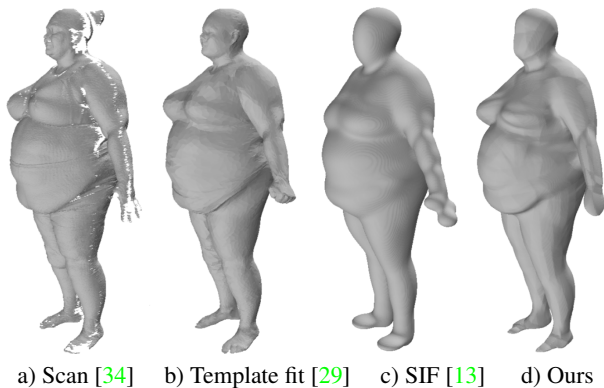


Figure 9. **Surface reconstruction from partial human scans.**

of local 3D functions, each evaluated as the product of a Gaussian and a residual function predicted with a deep network. We describe a method for inferring an LDIF from a 3D surface or posed depth image by first predicting a structured decomposition into shape elements, encoding 3D points within each shape element using PointNet [30], and decoding them with a small residual decoder. This approach provides an end-to-end framework for encoding shapes in local regions arranged in a global structure.

We show that this LDIF representation improves both reconstruction accuracy and generalization behavior over previous work – its F-Score results are better than the state-of-the-art [24] by 10.3 points for 3D autoencoding of test models from trained classes and by 17.8 points for unseen classes. We show that it dramatically reduces network parameter count – its local decoder requires approximately 0.4% of the parameters used by [24]. We show that it can be used to complete posed depth images – its depth completion results are 15.8 percentage points higher than [24]. Finally, we show that it can be used without change to reconstruct complete 3D surfaces of human bodies from partial scans.

Limitations and future work: Though the results are encouraging, there are limitations that require further investigation. First, we decompose space into a flat *set* of local regions – it would be better to consider a multiresolution hierarchy. Second, we leverage known camera poses when reconstructing shapes from depth images – it would be better to estimate them. Third, we estimate a constant number of local regions – it would be better to derive a variable number dynamically during inference (e.g., with an LSTM). Finally, we just scratch the surface of how structured and implicit representations can be combined – this is an interesting topic for future research.

Acknowledgements: We thank Boyang Deng for sharing OccNet* training code and Max Jiang for creating single-view depth renderers. We also thank Fangyin Wei and JP Lewis for feedback on the manuscript.

References

- [1] Brett Allen, Brian Curless, Brian Curless, and Zoran Popović. The space of human body shapes: reconstruction and parameterization from range scans. In *ACM transactions on graphics (TOG)*, volume 22, pages 587–594. ACM, 2003. [8](#)
- [2] Bruce G Baumgart. A polyhedron representation for computer vision. In *Proceedings of the May 19-22, 1975, national computer conference and exposition*, pages 589–596. ACM, 1975. [2](#)
- [3] James F Blinn. A generalization of algebraic surface drawing. *ACM transactions on graphics (TOG)*, 1(3):235–256, 1982. [2](#)
- [4] Jules Bloomenthal and Ken Shoemake. Convolution surfaces. *ACM SIGGRAPH Computer Graphics*, 25(4):251–256, 1991. [2](#)
- [5] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Generative and discriminative voxel modeling with convolutional neural networks. *NeurIPS 3D Deep Learning Workshop*, 2016. [2](#)
- [6] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing 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. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015. [4](#)
- [7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019. [1](#), [2](#)
- [8] Julian Chibane, Thiemo Alldieck, and Gerard Pons-Moll. Implicit functions in feature space for 3d shape reconstruction and completion. *arXiv preprint arXiv:2003.01456*, 2020. [2](#)
- [9] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. [2](#), [4](#)
- [10] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 605–613, 2017. [2](#)
- [11] James D Foley, Foley Dan Van, Andries Van Dam, Steven K Feiner, John F Hughes, J Hughes, and Edward Angel. *Computer graphics: principles and practice*, volume 12110. Addison-Wesley Professional, 1996. [2](#)
- [12] Lin Gao, Jie Yang, Tong Wu, Yu-Jie Yuan, Hongbo Fu, Yu-Kun Lai, and Hao Zhang. Sdm-net: Deep generative network for structured deformable mesh. *ACM Transactions on Graphics (TOG)*, 38(6):1–15, 2019. [2](#)
- [13] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7154–7164, 2019. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [14] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9785–9795, 2019. [2](#)
- [15] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 216–224, 2018. [2](#), [4](#), [5](#)
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778, 2016. [4](#), [11](#)
- [17] Chiyu Max Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, and Thomas Funkhouser. Local implicit grid representations for 3d scenes. *arXiv preprint arXiv:2003.08981*, 2020. [2](#)
- [18] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 371–386, 2018. [2](#)
- [19] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, pages 61–70. Eurographics Association, 2006. [2](#)
- [20] Jun Li, Kai Xu, Siddhartha Chaudhuri, Ersin Yumer, Hao Zhang, and Leonidas Guibas. Grass: Generative recursive autoencoders for shape structures. *ACM Transactions on Graphics (TOG)*, 36(4):52, 2017. [2](#)
- [21] Chen-Hsuan Lin, Chen Kong, and Simon Lucey. Learning efficient point cloud generation for dense 3d object reconstruction. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. [2](#)
- [22] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, oct 2015. [7](#)
- [23] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *Proceedings of the 14th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '87*, pages 163–169, New York, NY, USA, 1987. ACM. [4](#)
- [24] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019. [1](#), [2](#), [3](#), [4](#), [5](#), [8](#), [11](#)
- [25] Shigeru Muraki. Volumetric shape description of range data using blobby model. *ACM SIGGRAPH computer graphics*, 25(4):227–235, 1991. [2](#)
- [26] Yutaka Ohtake, Alexander Belyaev, Marc Alexa, Greg Turk, and Hans-Peter Seidel. Multi-level partition of unity implicits. In *Acm Siggraph 2005 Courses*, pages 173–es. 2005. [2](#)
- [27] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation.

- In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 1, 2
- [28] Despoina Paschalidou, Ali Osman Ulusoy, and Andreas Geiger. Superquadrics revisited: Learning 3d shape parsing beyond cuboids. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [29] Leonid Pishchulin, Stefanie Wuhrer, Thomas Helten, Christian Theobalt, and Bernt Schiele. Building statistical shape spaces for 3d human modeling. *Pattern Recognition*, 2017. 8
- [30] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 2, 4, 8, 11
- [31] Antonio Ricci. A constructive geometry for computer graphics. *The Computer Journal*, 16(2):157–160, 1973. 2
- [32] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3577–3586, 2017. 2
- [33] Gernot Riegler, Ali Osman Ulusoy, Horst Bischof, and Andreas Geiger. Octnetfusion: Learning depth fusion from data. In *2017 International Conference on 3D Vision (3DV)*, pages 57–66. IEEE, 2017. 2
- [34] Kathleen M Robinette, Sherri Blackwell, Hein Daanen, Mark Boehmer, and Scott Fleming. Civilian american and european surface anthropometry resource (caesar), final report. volume 1. summary. Technical report, SYTRONICS INC DAYTON OH, 2002. 8
- [35] Gopal Sharma, Rishabh Goyal, Difan Liu, Evangelos Kalogerakis, and Subhansu Maji. Csgnet: Neural shape parser for constructive solid geometry. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5515–5523, 2018. 2
- [36] Edward J Smith, Scott Fujimoto, Adriana Romero, and David Meger. Geometrics: Exploiting geometric structure for graph-encoded objects. *Proceedings of the 36th International Conference on Machine Learning*, 2019. 2
- [37] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2088–2096, 2017. 2
- [38] Maxim Tatarchenko, Stephan R Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3405–3414, 2019. 4
- [39] Shubham Tulsiani, Hao Su, Leonidas J Guibas, Alexei A Efros, and Jitendra Malik. Learning shape abstractions by assembling volumetric primitives. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2635–2643, 2017. 2
- [40] Gül Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning from synthetic humans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 6, 7
- [41] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 52–67, 2018. 2
- [42] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *Advances in neural information processing systems*, pages 82–90, 2016. 2
- [43] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015. 2
- [44] Geoff Wyvill, Craig McPheeters, and Brian Wyvill. Data structure for soft objects. In *Advanced Computer Graphics*, pages 113–128. Springer, 1986. 2
- [45] Qiangeng Xu, Weiyue Wang, Duygu Ceylan, Radomir Mech, and Ulrich Neumann. Disn: Deep implicit surface network for high-quality single-view 3d reconstruction. In *Advances in Neural Information Processing Systems*, pages 490–500, 2019. 1, 2
- [46] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4541–4550, 2019. 2

A. Hyperparameters

Table 5 contains all hyperparameter values used for training the model. Architecture details for individual networks are below.

ResNet50. We use a ResNet50 [16] V2 that is trained from scratch. The 20 depth images are concatenated channel-wise prior to encoding.

PointNet. We modify the original PointNet [30] architecture by removing the 64x64 orthogonal transformation to improve speed and reduce memory requirements.

OccNet. Our local decoder follows the same overall structure as the original OccNet [24]. However, we reduce the number of residual blocks from 5 to 1. The latent layer feature widths are also decreased proportionally to the vector dimensionality.

Local Point Cloud Extraction. We sample a subset of points for encoding by the local PointNet as follows. We first transform all 10,000 points to the local frame. Then we choose a distance threshold $r = 4.0$ measured in local units. Since the local frame is scaled proportionally to the radius, this threshold is approximately four radii in the world frame. We randomly sample 1,000 points without replacement within r and return those as the set of points to be encoded. If 1,000 points do not exist, we expand r until 1,000 total points are found.

Global Point Cloud Creation. In order to create 10,000 points from one or more input depth images, we randomly sample valid points without replacement from the depth images. If 10,000 valid pixels do not exist, we repeat random points as necessary before moving to the local extraction phase.

Activation Functions. Since the generated network activations \mathbf{y} are in the range $[-\infty, \infty]$, we apply activation functions to latents \mathbf{y} to interpret them as the analytic parameters θ_i . The following functions are used. For constants c_i : $-|y_{c,i}|$. For ellipsoid radii r_i : $0.15 \times \text{sig}(y_{r,i})$. For ellipsoid euler-angles e_i : $\max(\min(\frac{\pi}{4}, y_{e,i}), \frac{-\pi}{4})$. For ellipsoid positions p_i : $\frac{y_{p,i}}{2}$.

Metrics. Below we report details for how each metric is computed. All metrics are computed against the watertight version of the ground truth mesh in order to be consistent with the OccNet [24] procedure, and are computed in the normalized coordinate frame provided by [24]. Our predicted results are initially generated in a coordinate frame normalized using the centroid and variance of the mesh, rather than the bounding-box-based normalization of [24]. Therefore we transform back to the bounding-box normalized frame to compute metrics.

IoU. 100,000 uniform point samples with inside/outside labels are distributed by [24]. We evaluate the surface at these locations and compute the IoU between samples that are la-

Name	Value
α	100.0
w_S	0.1
w_U	1.0
w_C	10.0
w_P	1.0
Batch Size	24
Adam β_1	0.9
Adam β_2	0.999
Learning Rate	5×10^{-5}
Surface Isolevel	-0.07

Table 5. Hyperparameters and optimization details for training the autoencoder network.

beled inside by our representation and samples labeled as inside by the provided points.

F-Score. F-Score is computed at an absolute threshold of $\tau = 0.01$ units in the normalized coordinate space of [24]. We randomly sample 100,000 points on the surface of each mesh, and use point-to-point distances.

Chamfer. Chamfer distance is computed with the following formula. The factor of 100 is a scaling factor applied for consistency with existing approaches and readability.

$$100 * \left(\frac{1}{|A|} \sum_{a \in A} \min_{b \in B} \|a - b\|_2^2 + \frac{1}{|B|} \sum_{b \in B} \min_{a \in A} \|a - b\|_2^2 \right)$$