

Implicit Functions in Feature Space for 3D Shape Reconstruction and Completion

Julian Chibane^{1,2}Thiemo Alldieck^{1,3}Gerard Pons-Moll¹¹Max Planck Institute for Informatics, Saarland Informatics Campus, Germany²University of Würzburg, Germany³Computer Graphics Lab, TU Braunschweig, Germany

{jchibane, gpons}@mpi-inf.mpg.de alldieck@cg.cs.tu-bs.de

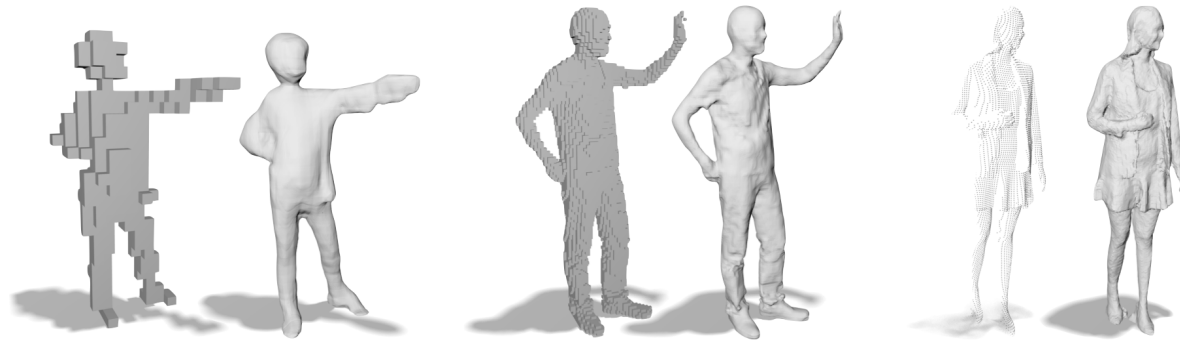


Figure 1: Results using our method. Left: sparse voxel reconstruction, middle: dense voxel reconstruction, right: 3D single-view point cloud reconstruction (back occluded). Our method delivers continuous outputs, handles multiple topologies (right) and unlike prior work, retains detail in the input (middle and right), and performs well with articulated humans.

Abstract

While many works focus on 3D reconstruction from images, in this paper, we focus on 3D shape reconstruction and completion from a variety of 3D inputs, which are deficient in some respect: low and high resolution voxels, sparse and dense point clouds, complete or incomplete. Processing of such 3D inputs is an increasingly important problem as they are the output of 3D scanners, which are becoming more accessible, and are the intermediate output of 3D computer vision algorithms. Recently, learned implicit functions have shown great promise as they produce continuous reconstructions. However, we identified two limitations in reconstruction from 3D inputs: 1) details present in the input data are not retained, and 2) poor reconstruction of articulated humans. To solve this, we propose Implicit Feature Networks (IF-Nets), which deliver continuous outputs, can handle multiple topologies, and complete shapes for missing or sparse input data retaining the nice properties of recent learned implicit functions, but critically they can also retain detail when it is present in the input data, and can reconstruct articulated humans. Our work differs from prior work in two crucial aspects. First, instead of using a single vector to encode a 3D shape, we extract a learnable 3-dimensional multi-scale tensor of deep fea-

tures, which is aligned with the original Euclidean space embedding the shape. Second, instead of classifying x - y - z point coordinates directly, we classify deep features extracted from the tensor at a continuous query point. We show that this forces our model to make decisions based on global and local shape structure, as opposed to point coordinates, which are arbitrary under Euclidean transformations. Experiments demonstrate that IF-Nets clearly outperform prior work in 3D object reconstruction in ShapeNet, and obtain significantly more accurate 3D human reconstructions. Code is available at <https://virtualhumans.mpi-inf.mpg.de/ifnets/>.

1. Introduction

While many works focus on image-based 3D reconstruction [23], in this paper, we focus on 3D surface reconstruction and shape completion from a variety of 3D inputs, which are deficient in some respect: low-resolution voxel-grids, high-resolution voxel-grids, sparse and dense point-clouds, complete or incomplete. Such inputs are becoming ubiquitous as 3D scanning technology is increasingly accessible, and they are often an intermediate output of 3D computer vision algorithms. However, the final output for most applications should be a renderable continuous

and complete surface, which is the focus of our work.

For sparse grids and (incomplete) point clouds, learning-based methods are a better choice than classical methods [6, 42], as they reason about global object shape, but are limited by their output representation. Mesh-based methods typically learn to deform an initial convex template [72], and hence can not represent different topologies. Voxel-based representations [11, 39] have a large memory footprint, which critically limits the output resolution to coarse shapes, without detail. Point cloud [54, 55] representations are more efficient but do not trivially enable rendering and visualization of the surfaces.

Recently, implicit functions [48, 43, 10] have shown to be a promising shape representation for learning. The key idea is to learn a function which, given a coarse shape encoded as a vector, and the x-y-z coordinates of a query point, decide whether the point is inside or outside of the shape. The learned implicit function can be evaluated at query 3D points at arbitrary resolutions, and the mesh/surface can be extracted applying the classical marching cubes algorithm. This output representation enables shape recovery at arbitrary resolutions, is continuous and can handle different topologies.

While these approaches work well to reconstruct aligned rigid objects, we observed they suffer from two main limitations: 1) they can not represent complex objects like articulated humans (reconstructions often miss arms or legs), 2) they do not retain detail present in the input data. We hypothesize this occurs because 1) networks learn an overly strong prior on x-y-z point coordinates damaging the invariance to articulation, and 2) the shape encoding vector lacks 3D structure, resulting in decodings that look more like classification into shape prototypes [66] rather than continuous regression. Consequently, all existing learning-based approaches, either based on voxels, meshes, points or implicit functions are lacking in some respect.

In this paper we propose *Implicit Feature Networks* (IF-Nets), which, unlike previous work, do well in 5 different axis, shown in Table 1: they are *continuous*, can handle *multiple topologies*, can complete data for *sparse input*, retaining the nice properties of implicit function models [48, 43, 10], but crucially they *also* retain detail when is present in the input (*dense input*), and can reconstructed *articulated* humans. IF-Nets *differ* from recent work [48, 10, 43] in two crucial aspects. First, instead of using a single vector to encode a 3D shape, we extract a 3-dimensional multi-scale tensor of deep features, which is aligned with original Euclidean space embedding the shape. Second, instead of classifying x-y-z point coordinates directly, we classify deep features extracted at continuous query points. Hence, unlike previous work, IF-nets do not *memorize* common x-y-z locations, which are arbitrary under Euclidean transformations. Instead, they make

Output 3D Repr.	Continuous Output	Multiple topologies	Sparse Input	Dense Input	Articulated
Voxels	✗	✓	✓	✗	✓
Points	✗	✓	✓	✗	✓
Meshes	✗	✗	✓	✗	✓
Implicit*	✓	✓	✓	✗	✗
Ours	✓	✓	✓	✓	✓

Table 1. Overview of strengths and weaknesses of recent 3D reconstruction approaches classified by their output representation. Voxels, point clouds, and meshes are non-continuous and suffer from discretization. Meshes additionally have fixed topologies, which limits the space of representable 3D shapes. Recent learned implicit functions* [43, 10, 48] alleviate these limitations but fail to retain details or reconstruct articulation. The proposed IF-Nets share the desired properties of implicit functions for reconstructing from 3D inputs, but are *additionally* able to preserve detail present in dense 3D input and to reconstruct articulated humans.

decisions based on multi-scale features encoding local and global object shape structures around the point.

To demonstrate the advantages of IF-Nets, first, we show that IF-Nets can reconstruct simple rigid 3D objects at better accuracy than previous methods. In ShapeNet [9], IF-Nets outperform the state-of-the-art results. For articulated humans, we train IF-Nets and related methods on a dataset of 1600 humans in varied poses, shapes, and clothing. In stark contrast to recent work [48, 10, 43], IF-Nets can reconstruct articulated objects globally without missing limbs, while recovering detailed structures such as cloth wrinkles. Quantitative and qualitative experiments validate that IF-Nets are more robust to articulations and produce globally consistent shapes without losing fine-scale detail. To encourage further research in 3D processing, learning, and reconstruction, we make IF-Nets publicly available at <https://virtualhumans.mpi-inf.mpg.de/ifnets/>.

2. Related Work

Approaches for *3D shape reconstruction* can be classified according to the representation used: voxels, meshes, point clouds, and implicit functions; and according to the object type: rigid objects vs humans. For a more exhaustive recent review, we refer the reader to [23]. A condensed overview of strengths and weaknesses of recent 3D reconstruction approaches is given in Table 1.

Voxels for rigid objects: Since voxels are a natural 3D extension to pixels in image grids and admit 3D convolutions, they are most commonly used for generation and reconstruction [29, 26, 57, 46]. However, the memory footprint scales cubically with the resolution, which limited early works [75, 11, 68] to predict shapes in small 32^3 grids. Higher resolutions have been used [74, 73, 81] at the cost of limited training batches and slow training or lossy 2D projections [63]. Multi-resolution [24, 65, 71] reconstruction reduced the memory footprint, allowing grids of size

256^3 . However, the approaches are complicated to implement, require multiple passes over the input, and are still limited to grids of size 256^3 , which result in visible quantization artifacts. To smooth out noise, it is possible to represent shapes as Truncated Signed Distance functions [12] for learning [14, 36, 58, 64]. The resolution is however still bounded by the 3D grid storing the TSDF values.

Generative shape models typically map a 1D vector to a voxel representation with a neural network [18, 74]. Like us, the authors of [40] observe that the 1D vector is too restrictive to generate shapes with global and local structures. They introduce a hierarchical latent code with skip connections. Instead, we propose a much simpler 3-dimensional multi-scale feature tensor, which is aligned with original Euclidean space embedding the shape.

Humans with voxels: From images, CNN based reconstruction of humans represented as voxels [70, 17, 82] or depth-maps [16, 62, 38] typically produce more details than mesh or template-based representations, because predictions are aligned with the input pixels. Unfortunately, this comes at the cost of missing parts in the body. Hence, some methods [70, 62] fit the SMPL [41] model to the reconstructions as a post-processing step. This is however prone to fail if the original reconstructions are too incomplete. All these approaches process image pixels whereas we focus on processing 3D data directly. Unlike our IF-Nets, these methods are bounded by the resolution of the voxel grid.

Meshes for rigid objects: Most mesh-based methods predict shape as a deformation from a template [72, 56] and hence are limited to a single topology. Alternatively, the mesh (vertices and faces) can be inferred directly [20, 13] – while this research direction is promising, methods are still computationally expensive and can not guarantee a closed mesh without intersections. Direct mesh prediction can also be obtained using a learnable version [40] of the classical marching cubes algorithm [42], but the approach is limited to an underlying small voxel grid of 32^3 . Promising combinations of voxels and meshes have been proposed [19], but results are still coarse.

Meshes for humans: Since the introduction of the (mesh-based) SMPL human model [41] there have been a growing number of papers leveraging it to reconstructing shape and pose from point clouds, depth data and images [28, 32, 33, 47, 69, 79]. Since SMPL does not model clothing and detail, recent methods predict deformations from SMPL [1, 2, 3, 51, 7] or a template [21, 22]. Unfortunately, CNN based mesh predictions tend to be over-smooth. More detail can be obtained predicting normals and displacement maps on a UV-map/geometry image of the surface [4, 37, 53]. However, all these approaches require different templates [7, 49] for every new garment topology or do not produce high-quality reconstructions [53].

Point clouds for rigid objects: Processing point clouds is an important problem as they are the output of many sensors (LiDAR, 3D scanners) and computer vision algorithms. Due to their low weight, they have been also popular in computer graphics for representing and manipulating shapes [50]. PointNet based architectures [54, 55] were the first to process point clouds directly for classification and semantic segmentation. The idea is to apply a fully connected network to each point followed by a global pooling operation to achieve permutation invariance. Recent architectures apply kernel point convolutions [67], tree-based graph convolutions [60], and normalizing flows [77]. Point clouds are also used as shape representation for reconstruction [15, 25] and generation [77]. Unlike voxels or meshes, point clouds need to be non-trivially post-processed using classical methods [6, 30, 31, 8] to obtain renderable surfaces.

Point clouds for humans: Very few works represent humans with point clouds [5], probably because they can not be rendered. Recent works have employed either PointNet architectures [27] or architectures based on point bases [52] to register a human mesh to the point cloud.

Implicit Functions for rigid objects: Recently, neural networks have been used to learn a continuous implicit function representing shape [43, 48, 10, 44, 35]. For this a neural network can be feed with a latent code and a query point (x - y - z) to predict the TSDF value [48] or the binary occupancy of the point [43, 10]. A recent method [76] achieved state-of-the-art results for 3D reconstruction from images combining 3D query point features with local image features, by approximating the projection of the query point onto the 2D image with a view-point prediction. This trick of querying continuous points used in implicit function learning allows predicting in continuous space (potentially at any resolution), breaking the memory barrier of voxel-based methods. These works inspired our work, but we note that they can not reconstruct articulated humans from 3D data: [76] can not take 3D inputs as point clouds or voxel grids and relies on an approximate 3D to 2D projection losing details; the reconstructions of [43, 10] often miss limbs. We hypothesize that [43, 10, 76] memorize point coordinates instead of reasoning about shape, and that the vectorized latent 1D vector representation [43, 10] is not aligned with the input, and lacks 3D structure. We address this issues by querying deep features extracted at continuous locations of a 3D grid of multi-scale features aligned with the 3D input space. This modification is easy to implement and results in significant gains in reconstruction quality.

Implicit Functions for humans: TSDFs [12] have been used to represent human shapes for depth-fusion and tracking [45, 61]. Such implicit representation has been combined with the SMPL [41] body model to significantly increase tracking robustness and accuracy [78]. From an in-

put image, humans in clothing are predicted using an implicit network [59]. [59] produces higher quality results compared to prior implicit function work [10]. The reconstruction is done by pointwise occupancy prediction based on the location of a 3D query point and 2D image features. For simple poses, the approach produces very compelling and detailed results but struggles for more complex poses. The approach [59] does not incorporate a multi-scale 3D shape representation like ours, and it is designed for image reconstruction, whereas we focus on 3D reconstruction from sparse and dense point clouds and occupancy grids. Like previous implicit networks, our method produces continuous surfaces at arbitrary resolution. But importantly, by virtue of our 3D multi-scale shape representation aligned with the input space, our reconstructions preserve global structure while retaining fine-scale detail, even for complex poses.

3. Method

To motivate the design of our Implicit Feature Networks (IF-Nets), we first describe the formulation of recent learned implicit functions, pointing out their strengths and weaknesses in Sec. 3.1. We explain our IF-Nets in Sec. 3.2. The key ideas of IF-Nets are illustrated in Fig. 2.

3.1. Background: Learning with Implicit Surfaces

While recent works [48, 43, 10] on learned implicit reconstruction from 3D input differ in their inference and their output shape representation (signed distance or binary occupancies), they are conceptually similar. Here, we describe the occupancy formulation of [43]. Note that the strengths and limitations of these methods are very similar. They all encode a 3D shape using a latent vector $\mathbf{z} \in \mathcal{Z} \subset \mathbb{R}^m$. Then a continuous representation of the shape is obtained by learning a neural function

$$f(\mathbf{z}, \mathbf{p}): \mathcal{Z} \times \mathbb{R}^3 \mapsto [0, 1], \quad (1)$$

which given a query point $\mathbf{p} \in \mathbb{R}^3$, and the latent code \mathbf{z} , classifies whether the point is inside (classification as 1) or outside (classification as 0) the surface. Thereby, the surface is implicitly represented as the points on the decision boundary, $\{\mathbf{p} \in \mathbb{R}^3 \mid f(\mathbf{z}, \mathbf{p}) = t\}$, with a threshold parameter t ($t = 0.5$ for IF-Nets).

Once $f(\cdot)$ is learned, it can be queried at continuous point locations, without resolution restrictions imposed by typical voxel grids. To construct a mesh, marching cubes [42] can be applied on the predicted occupancy grid. This elegant formulation breaks the barriers of previous representations allowing detailed reconstruction of complex topologies, and has proven effective for several tasks such as rigid object reconstruction from images, occupancy grids, and point clouds. However, we observed that models

of this kind suffer from two main limitations: 1) they can not represent complex objects like articulated objects, and 2) they do not preserve detail present in the input data. We address these limitations with IF-Nets.

3.2. Implicit Feature Networks

We identify two potential problems with the previous formulation. First, directly inputting point coordinates \mathbf{p} gives the network the option to by-pass reasoning about shape structure, by memorizing typical point occupancies for object prototypes. This severely damages reconstruction in-variance to rotation and translation, which is one of the cornerstones of successful 2D convolution networks for segmentation, recognition, and detection. Second, encoding the full shape in a single vector \mathbf{z} loses detail present in the data, and loses alignment with the original 3D space where shapes are embedded.

In this work, we propose a novel encoding and decoding tandem capable of addressing the above limitations for the task of 3D reconstruction from point clouds or occupancy grids. Given such 3D input data $\mathbf{X} \in \mathcal{X}$ of an object, where \mathcal{X} denotes the space of the inputs, and a 3D point $\mathbf{p} \in \mathbb{R}^3$, we want to predict if \mathbf{p} lies inside or outside the object.

Shape Encoding: Instead of encoding the shape in a single vector \mathbf{z} , we construct a rich encoding of the data \mathbf{X} through subsequently convolving it with learned 3D convolutions. This requires the input to lie on a discrete voxel grid, i.e. $\mathcal{X} = \mathbb{R}^{N \times N \times N}$, where $N \in \mathbb{N}$ denotes the input resolution. To process point clouds we simply discretize them first. The convolutions are followed by down scaling the input, creating growing receptive fields and channels but shrinking resolution, just like commonly done in 2D [34]. Applying this procedure recursively n times on the input data \mathbf{X} , we create *multi-scale deep feature grids* $\mathbf{F}_1, \dots, \mathbf{F}_n$, $\mathbf{F}_k \in \mathcal{F}_k^{K \times K \times K}$, of decreasing resolution $K = \frac{N}{2^{k-1}}$, and variable channel dimensionality $F_k \in \mathbb{N}$ at each stage $\mathcal{F}_k \subset \mathbb{R}^{F_k}$. The feature grids \mathbf{F}_k at the early stages (starting at $k = 1$) capture high frequencies (shape detail), whereas feature grids \mathbf{F}_k at the late stages (ending at stage $k = n$) have a large receptive fields, which capture the global structure of the data. This enables to reason about missing or sparse data, while retaining detail when is present in the input. We denote the encoder as

$$g(\mathbf{X}) := \mathbf{F}_1, \dots, \mathbf{F}_n. \quad (2)$$

Shape Decoding: Instead of classifying point coordinates \mathbf{p} directly, we extract the learned deep features $\mathbf{F}_1(\mathbf{p}), \dots, \mathbf{F}_n(\mathbf{p})$ from the feature grids at location \mathbf{p} . This is only possible because our *encoding has a 3D structure* aligned with the input data. Since feature grids are discrete, we use trilinear interpolation to query continuous 3D points

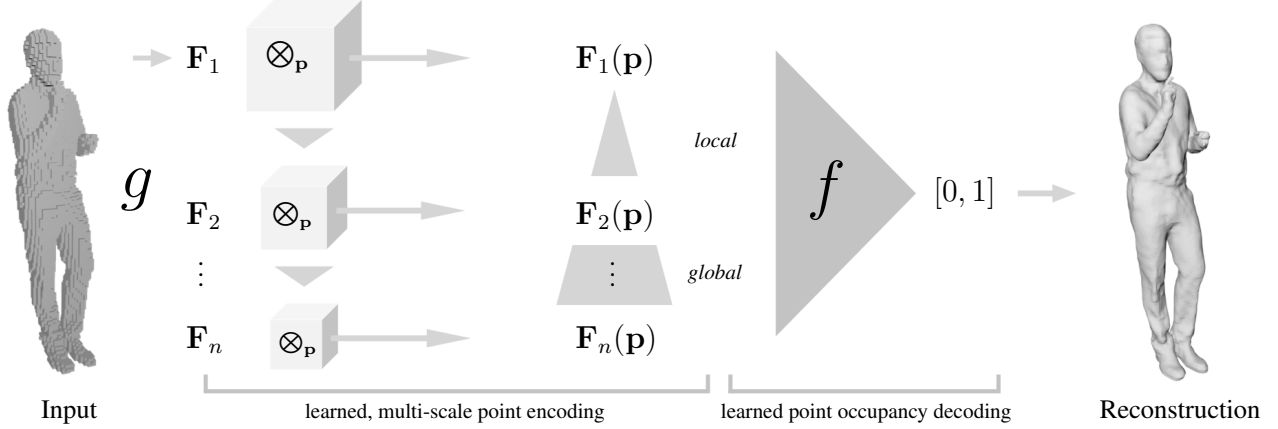


Figure 2. Overview of IF-Nets: given an (incomplete or low resolution) input, we compute a 3D grid of multi-scale features, encoding global and local properties of the input shape. Then, we extract deep features $\mathbf{F}_1(\mathbf{p}) \dots \mathbf{F}_n(\mathbf{p})$ from the grid at *continuous point* locations \mathbf{p} . Based *only on these features* a decoder $f(\cdot)$ decides whether the point \mathbf{p} lies inside (classification as 1) or outside (classification as 0) the surface. Like recent implicit function-based works, we can query at arbitrary resolutions and reconstruct a continuous surface. Unlike them, our method reasons based exclusively on point-wise deep features, instead of point coordinates. This allows us to reconstruct articulated structures and preserve input detail.

$\mathbf{p} \in \mathbb{R}^3$. In order to encode information of the local neighborhood into the point encoding, even at early grids with small receptive fields (e.g. \mathbf{F}_1), we extract features at the location of a query point \mathbf{p} itself and *additionally* at surrounding points in a distance d along the Cartesian axes:

$$\{\mathbf{p} + a \cdot \mathbf{e}_i \cdot d \in \mathbb{R}^3 \mid a \in \{1, 0, -1\}, i \in \{1, 2, 3\}\}, \quad (3)$$

where $d \in \mathbb{R}$ is the distance to the center point \mathbf{p} and $\mathbf{e}_i \in \mathbb{R}^3$ is the i -th Cartesian axis unit vector, see supplementary material for an illustration.

The point encoding $\mathbf{F}_1(\mathbf{p}), \dots, \mathbf{F}_n(\mathbf{p})$, with $\mathbf{F}_k(\mathbf{p}) \in \mathcal{F}_k$, is then fed into a point-wise decoder $f(\cdot)$, parameterized by a fully connected neural network, to predict if the point \mathbf{p} lies inside or outside the shape:

$$f(\mathbf{F}_1(\mathbf{p}), \dots, \mathbf{F}_n(\mathbf{p})): \mathcal{F}_1 \times \dots \times \mathcal{F}_n \mapsto [0, 1] \quad (4)$$

In contrast to Eq. (1), in this formulation, the network classifies the point based on local and global shape features, instead of point coordinates, which are arbitrary under rotation, translation, and articulation transformations. Furthermore, due to our multi-scale encoding, details can be preserved while reasoning about global shape is still possible.

3.3. Method Training

To train the multi-scale encoder $g_{\mathbf{w}}(\cdot)$ in Eq. (2), and decoder $f_{\mathbf{w}}(\cdot)$ in Eq. (4), parameterized with neural weights \mathbf{w} , pairs $\{\mathbf{X}_i, \mathcal{S}_i\}_{i=1}^T$ of 3D inputs \mathbf{X}_i with corresponding 3D ground truth object surfaces \mathcal{S}_i are required, where $i \in 1, \dots, T$ and $T \in \mathbb{N}$ denotes the number of such training examples. The notation $g_{\mathbf{w}}(\mathbf{X}, \mathbf{p}) := \mathbf{F}_1^{\mathbf{w}}(\mathbf{p}), \dots, \mathbf{F}_n^{\mathbf{w}}(\mathbf{p})$ denotes evaluation of the multi-scale encoding at point \mathbf{p} .

To create training point samples, for every ground truth surface \mathcal{S}_i , we sample a number $S \in \mathbb{N}$ of points $\mathbf{p}_i^j \in \mathbb{R}^3$, $j \in 1, \dots, S$. To this end, we first make the ground truth surface \mathcal{S}_i watertight. Then we compute the ground truth occupancy $o_i(\mathbf{p}_i^j) \in \{0, 1\}$, which evaluates to 1 for inside points and 0 otherwise. Next, the point samples \mathbf{p}_i^j are created near the surface by sampling points $\mathbf{p}_{i,j}^S \in \mathcal{S}_i$ on the ground truth surfaces and adding random displacements $\mathbf{n}_{i,j} \sim \mathcal{N}(0, \Sigma)$, i.e. $\mathbf{p}_i^j := \mathbf{p}_{i,j}^S + \mathbf{n}_{i,j}$. To this end, we use a diagonal covariance matrix $\Sigma \in \mathbb{R}^{3 \times 3}$ with entries $\Sigma_{i,i} = \sigma$. We find good results by sampling 50% of the point samples very near the surface with a small σ_1 , and 50% in the further away surroundings with a larger σ_2 . For training, the network weights \mathbf{w} are optimized by minimizing the mini-batch loss

$$\begin{aligned} \mathcal{L}_{\mathcal{B}}(\mathbf{w}) &:= \sum_{i \in \mathcal{B}} \sum_{j \in \mathcal{R}} L(f_{\mathbf{w}}(g_{\mathbf{w}}(\mathbf{X}_i, \mathbf{p}_i^j)), o_i(\mathbf{p}_i^j)) \\ &= \sum_{i \in \mathcal{B}} \sum_{j \in \mathcal{R}} L(f_{\mathbf{w}}(\mathbf{F}_1^{\mathbf{w}}(\mathbf{p}_i^j), \dots, \mathbf{F}_n^{\mathbf{w}}(\mathbf{p}_i^j)), o_i(\mathbf{p}_i^j)), \end{aligned} \quad (5)$$

which sums over training surfaces $i \in \mathcal{B} \subset 1, \dots, T$ of a given mini-batch \mathcal{B} and point samples $j \in \mathcal{R} \subset 1, \dots, S$ of a subsample \mathcal{R} . The subsample \mathcal{R} is regenerated for every evaluation of the mini-batch loss $\mathcal{L}_{\mathcal{B}}$. For $L(\cdot, \cdot)$, we use the standard cross-entropy loss. By minimizing $\mathcal{L}_{\mathcal{B}}$, we train the encoder $g_{\mathbf{w}}(\cdot)$ and the decoder $f_{\mathbf{w}}(\cdot)$ jointly and end-to-end. Please see the supplementary material for the concrete values for the hyperparameters used in the experiments.

3.4. Method Inference

At test time, the goal is to reconstruct a continuous and complete representation, given only a discrete and incom-

plete 3D input \mathbf{X} . First, we use the learned encoder network to construct the multi-scale feature grids $g(\mathbf{X}) = \mathbf{F}_1, \dots, \mathbf{F}_n$. Then, we use the point-wise decoder network $f(g(\mathbf{X}, \mathbf{p}))$ to create occupancy predictions at continuous point locations $\mathbf{p} \in \mathbb{R}^3$ (cf. Sec. 3.2). In order to construct a mesh, we evaluate the IF-Net on points on a grid of the desired resolution. Then, the resulting high resolution occupancy grid is transformed into a mesh using the classical marching cubes [42] algorithm.

4. Experiments

In this section we validate the effectiveness of IF-Nets on the challenging task of 3D shape reconstruction. We show that our IF-Nets are able to address two limitations of recent learning-based approaches for this task: 1) IF-Nets preserve detail present in the input data, while also reasoning about incomplete data, 2) IF-Nets are able to reconstruct articulated humans in complex clothing. To this end, we conduct three experiments of increasing complexity: *Point Cloud Completion* (Sec. 4.1), *Voxel Super-Resolution* (Sec. 4.2) and *Single-View Human Reconstruction* (Sec. 4.3).

Baselines: For the task of Point Cloud Completion, we evaluate our approach against Occupancy Networks [43] (OccNet), Point Set Generation Networks [15] (PSGN) and Deep Marching Cubes [39] (DMC). For Voxel Super-Resolution, we compare against IMNET [10] as well as again against OccNet and DMC. For DMC and PSGN we used the implementations provided online by the authors of [43]. We trained all methods until the validation minimum was reached. Training was repeated for every considered experiment setup. To show a consistent comparison, we modified the IMNET implementation to be able to be trained on all ShapeNet classes jointly. For IMNET and OccNet, we kept the sampling strategies proposed by their authors. For IMNET, we followed the authors and performed progressive resolution increasing of training data sampling during training.

Metrics: To measure reconstruction quality quantitatively, we consider three established metrics (see suppl. of [43] for definition and implementation): volumetric *intersection over union* (IoU) measuring how well the defined volumes match (higher is better), *Chamfer- L_2* measuring the accuracy and completeness of the surface (lower is better), and *normal consistency* measuring the accuracy and completeness of the shape normals (higher is better).

Data: We consider two datasets: 1) a dataset containing 3D scans of humans¹ to assess the challenging task of reconstruction from incomplete and articulated shapes and 2) the established ShapeNet [9] dataset, consisting of rigid object classes, with rather prototypical shapes like cars, airplanes, and rifles. The ShapeNet data has been pre-

¹The dataset will be available for purchase from Twindom.

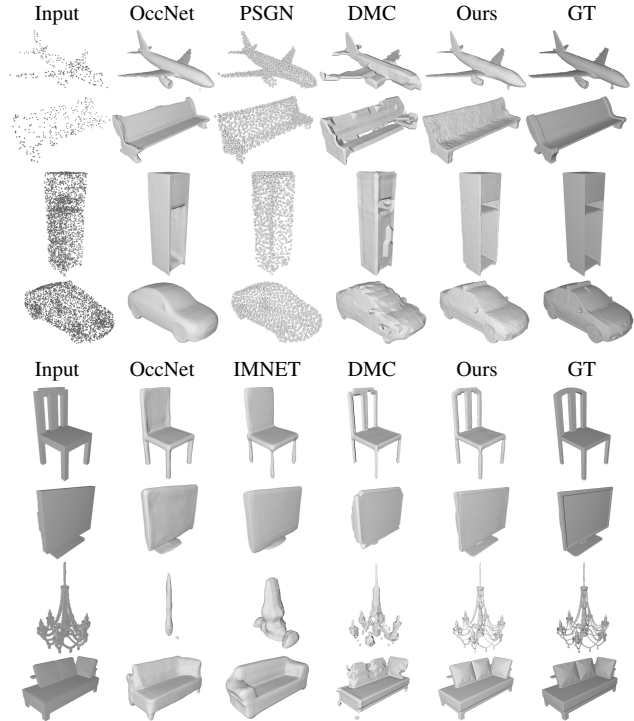


Figure 3. Qualitative results for two input types: point cloud (top) and voxels (bottom) on ShapeNet dataset. Each type is further subdivided into sparse (top two rows) and dense (bottom two rows).

processed to be watertight by the authors of [76], allowing to compute ground truth occupancies and scaled such that every shape’s largest bounding box edge has length one. We conduct all experiments and evaluations using pre-processed ShapeNet data and use the common training and test split by [11]. However, preprocessing failed for some objects, leading to broken objects with large holes. Therefore, 508 heavily distorted objects have been removed for meaningful evaluation. The filtered list of all used objects is published alongside the code. We also evaluate on a challenging dataset consisting of scanned humans in highly varying articulations with complex and varying clothing topologies like coats, skirts, or hats. The scans have been captured using commercial 3D scanners. The dataset, referred to as *Humans*, consists of 2183 such scans, split into 478 examples for testing, 1598 for training and 197 for validation. The scans have been height normalized and centered, but in contrast to the ShapeNet objects, exhibit varying rotations.

4.1. Point Cloud Completion

As a first task, we apply IF-Nets to the problem of completing sparse and dense point clouds – we sample 300 points (sparse) and 3000 points (dense) respectively from ShapeNet surface models and ask our method to complete the full surfaces. Completing point clouds is challenging

	IoU \uparrow		Chamfer- L_2 \downarrow		Normal-Consis. \uparrow	
Input	—	—	0.07	0.009	—	—
OccNet	0.73	0.72	0.03	0.04	0.88	0.88
DMC	0.58	0.65	0.03	0.01	0.83	0.86
PSGN	—	—	0.04	0.04	—	—
Ours	0.79	0.88	0.02	0.002	0.90	0.95

Table 2. Results of point cloud reconstruction on ShapeNet. Left number indicates score from 300 points, right one from 3000 points. Chamfer- L_2 results $\times 10^{-2}$.

	IoU \uparrow		Chamfer- L_2 \downarrow		Normal-Consis. \uparrow	
Input	0.49	0.79	0.04	0.003	0.81	0.87
DMC	0.59	0.67	0.45	0.45	0.83	0.84
IMNET	0.49	0.40	0.47	0.40	0.79	0.77
OccNet	0.60	0.71	0.10	0.05	0.85	0.88
Ours	0.73	0.92	0.02	0.002	0.91	0.98

Table 3. Results of voxel grid reconstruction on ShapeNet. For each metric, left column indicates score from 32^3 resolution, right one from 128^3 resolution. Chamfer- L_2 results $\times 10^{-2}$.

	IoU \uparrow		Chamfer- L_2 \downarrow		Normal-Consis. \uparrow	
Input	0.49	0.76	0.04	0.003	0.82	0.86
DMC	0.77	0.85	0.03	0.01	0.79	0.83
IMNET	0.63	0.64	0.27	0.23	0.79	0.79
OccNet	0.63	0.65	0.22	0.19	0.79	0.79
Ours	0.80	0.96	0.02	0.001	0.86	0.94

Table 4. Results of voxel grid reconstruction on the Humans dataset. Left number indicates score from 32^3 resolution, right one from 128^3 resolution. Chamfer- L_2 results $\times 10^{-2}$. IF-Nets coherently outperform others in the incomplete data setup. IF-Nets show a large increase in performance with dense data, whereas others show similar performance. This demonstrates that IF-Nets are the first learned approach, to our knowledge, being able to faithfully reconstruct dense information present in 3D data.

since it requires to simultaneously preserve input details and reason about missing structure at the same time. In Fig. 3 we show comparisons against the baseline methods. Our method outperforms all baselines both in preserving local detail and recovering global structures. For the dense point clouds, the strengths of our method are paramount. Our method is the only one capable of reconstructing the car rear-view mirrors and the additional shelf of the wardrobe. We additionally quantitatively compare our method and report the numbers in Tab. 2. Our method beats the state-of-the-art in all metrics by large margin. In fact, using 3000 points as input, all competitors produce results which have larger Chamfer distance than the input itself, suggesting they fail at preserving input detail. Only IF-Nets preserve input details while completing missing structures.

4.2. Voxel Super-Resolution

As a second task, we apply our method to 3D super-resolution. To effectively solve this task, our method needs to again preserve the input shape while reconstructing details not present in the input. Our results in side-by-side comparison with the baselines are depicted in Fig. 3 (bot-

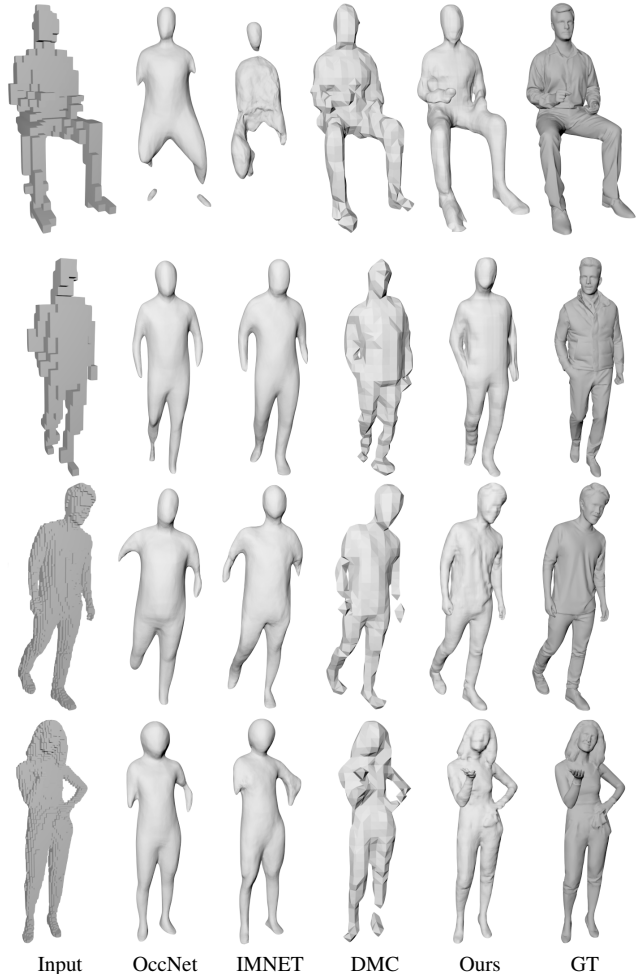


Figure 4. Qualitative results of sparse (32^3 , upper) and dense (128^3 , lower) 3D voxel super-resolution on the Humans dataset.

tom). While most baseline methods either hallucinate structure or completely fail, our method consistently produces accurate and highly detailed results. This is also reflected in the numerical comparison in Tab. 3, where we improve over the baselines in all metrics.

The two last examples in Fig. 3 illustrate the limitations of current implicit methods: If a shape differs too much from the training set, the method fails or seems to return a similar previously seen example. Consequentially, we hypothesize that the current methods are not suited for tasks where classification into shape prototypes is not sufficient. This is for example the case for humans as they come in various shapes and articulations. To verify our hypothesis, we additionally perform 3D super-resolution on our Humans dataset. Here the advantages are even more prominent: Our method is the only one that consistently reconstructs all limbs and produces highly detailed results. Implicit learning-based baselines produce truncated or completely missing limbs. We outperform all baselines also quantitatively (see Tab 4).

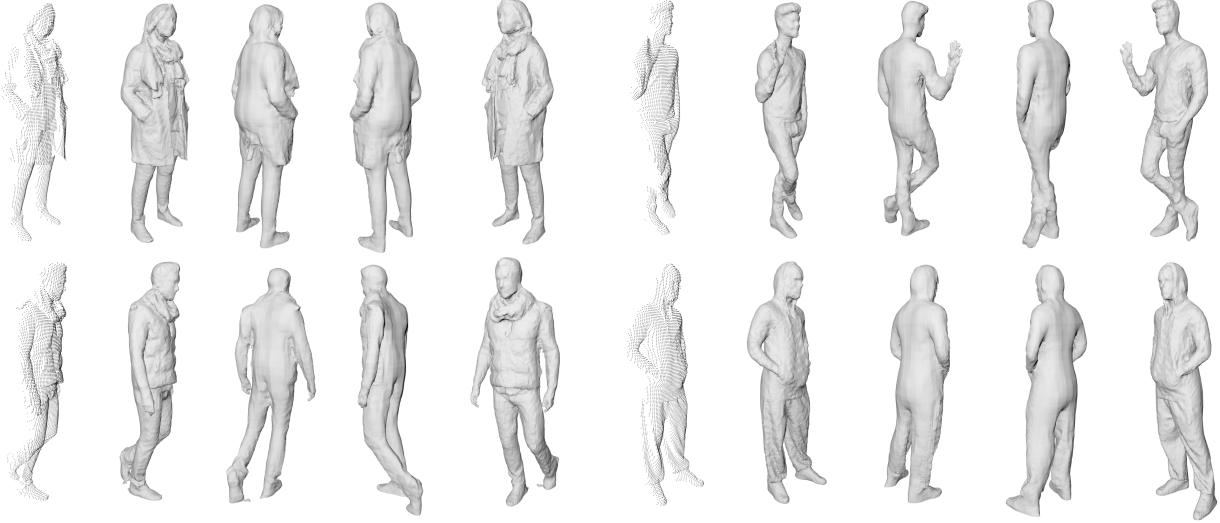


Figure 5. 3D single view reconstructions from point clouds (note that the back is completely occluded). For four different single-view point clouds, we show our reconstructions from four different viewpoints.

4.3. Single-View Human Reconstruction

Finally, to demonstrate the full capabilities of IF-Nets, we use them for single-view human reconstruction. In this task, only a partial 3D point cloud is given as input – the typical output of a depth camera. We conduct this experiment on the challenging Humans dataset, by rendering a 250×250 resolution depth image, yielding around 5000 points on the visible side of the subject. To successfully fulfill this task, our model has to simultaneously reconstruct novel articulations, retain fine details, and complete the missing data at the occluded regions – the input contains only one side of the underlying shape. Despite these challenges, our model is capable of reconstructing plausible and highly detailed shapes. In Fig. 5, we show both input and our results from four different angles. Note how fine structures like the scarfs, wrinkles, or individual fingers are present in the reconstructed shapes. Although the backside region (occluded) has less details than the visible one, IF-Nets always produce plausible surfaces.

This can also be seen quantitatively. *Ours*: IoU 0.86, Chamfer- L_2 0.011×10^{-2} , Normal-Consistency 0.90. *Input point cloud*: Chamfer- L_2 0.252×10^{-2} . The quantitative results are in between the reconstruction quality of 32^3 and 128^3 full subject voxel inputs (see Tab. 4), which once more validates that IF-Nets can complete single-view data. In the supplementary video, we show an additional result on single-view reconstruction on the BUFF dataset [80] from video (without retraining nor fine tuning the model).

5. Discussion and Conclusion

In this work, we have introduced IF-Nets for 3D reconstruction and completion from deficient 3D inputs. First, we

have argued for an encoding consisting of a 3D multi-scale tensor of deep features, which is aligned with the Euclidean space embedding the shape. Second, instead of classifying x-y-z coordinates directly, we classify deep features extracted at their location. Experiments demonstrate that IF-Nets deliver continuous outputs, can reconstruct multiple topologies such as 3D humans in varied clothing, and 3D objects from ShapeNet. Quantitatively, IF-Nets outperform all state-of-the-art baselines by a large margin in all tasks. Our reconstruction from single-view point clouds (detailed on the visible part but with missing data on the occluded part), demonstrate the strengths of IF-Nets: details in the input are preserved, while the shape is *completed* on the occluded part, even for articulated shapes.

Future work will explore extending IF-Nets to be generative, that is being able to sample detailed hypothesis conditioned on partial input. We also plan to address image-based reconstruction in 2 stages: first predicting a depth map, and then completing shape with IF-Nets.

With a rising number of computer vision image reconstruction methods producing partial 3D point clouds and voxels, and 3D scanners and depth cameras becoming accessible, 3D (deficient and incomplete) data will be omnipresent in the future, and IF-Nets have the potential to be an important building block for its reconstruction and completion.

Acknowledgments. We would like to thank Verica Lazova for helping creating the figures, Bharat Lal Bhatnagar for helping with data preprocessing, and Lars Mescheder for sharing their watertight Shapenet meshes. This work is funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - 409792180 (Emmy Noether Programme, project: Real Virtual Humans). We would like to thank Twindom for providing us with the scan data.

References

- [1] Thiemo Alldieck, Marcus Magnor, Bharat Lal Bhatnagar, Christian Theobalt, and Gerard Pons-Moll. Learning to reconstruct people in clothing from a single RGB camera. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 3
- [2] Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Detailed human avatars from monocular video. In *International Conference on 3D Vision*, 2018. 3
- [3] Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Video based reconstruction of 3D people models. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 3
- [4] Thiemo Alldieck, Gerard Pons-Moll, Christian Theobalt, and Marcus Magnor. Tex2shape: Detailed full human body geometry from a single image. In *IEEE International Conference on Computer Vision*. IEEE, 2019. 3
- [5] Tristan Aumentado-Armstrong, Stavros Tsogkas, Allan Jepson, and Sven Dickinson. Geometric disentanglement for generative latent shape models. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8181–8190, 2019. 3
- [6] Fausto Bernardini, Joshua Mittleman, Holly E. Rushmeier, Cláudio T. Silva, and Gabriel Taubin. The ball-pivoting algorithm for surface reconstruction. *IEEE Transactions on Visualization and Computer Graphics*, 5(4):349–359, 1999. 2, 3
- [7] Bharat Lal Bhatnagar, Garvita Tiwari, Christian Theobalt, and Gerard Pons-Moll. Multi-garment net: Learning to dress 3d people from images. In *IEEE International Conference on Computer Vision*. IEEE, 2019. 3
- [8] Fatih Calakli and Gabriel Taubin. SSD: smooth signed distance surface reconstruction. *Computer Graphics Forum*, 30(7):1993–2002, 2011. 3
- [9] 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. *arXiv preprint arXiv:1512.03012*, 2015. 2, 6
- [10] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019. 2, 3, 4, 6
- [11] 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 *European Conference on Computer Vision*, pages 628–644, 2016. 2, 6
- [12] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Computer Graphics and Interactive Techniques*, pages 303–312, 1996. 3
- [13] Angela Dai and Matthias Nießner. Scan2mesh: From unstructured range scans to 3d meshes. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5574–5583, 2019. 3
- [14] Angela Dai, Charles Ruizhongtai Qi, and Matthias Nießner. Shape completion using 3d-encoder-predictor cnns and shape synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 6545–6554, 2017. 3
- [15] Haoqiang Fan, Hao Su, and Leonidas J. Guibas. A point set generation network for 3d object reconstruction from a single image. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2463–2471, 2017. 3, 6
- [16] Valentin Gabeur, Jean-Sébastien Franco, Xavier Martin, Cordelia Schmid, and Grégory Rogez. Moulding humans: Non-parametric 3d human shape estimation from single images. In *IEEE International Conference on Computer Vision*, pages 2232–2241, 2019. 3
- [17] Andrew Gilbert, Marco Volino, John P. Collomosse, and Adrian Hilton. Volumetric performance capture from minimal camera viewpoints. In *European Conference on Computer Vision*, pages 591–607, 2018. 3
- [18] Rohit Girdhar, David F. Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and generative vector representation for objects. In *European Conference on Computer Vision*, pages 484–499, 2016. 3
- [19] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh R-CNN. In *IEEE International Conference on Computer Vision*, pages 9785–9795, 2019. 3
- [20] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 216–224, 2018. 3
- [21] Marc Habermann, Weipeng Xu, , Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Livecap: Real-time human performance capture from monocular video. *ACM Transactions on Graphics, (Proc. SIGGRAPH)*, jul 2019. 3
- [22] Marc Habermann, Weipeng Xu, , Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Deepcap: Monocular human performance capture using weak supervision. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2020. 3
- [23] Xian-Feng Han, Hamid Laga, and Mohammed Bennamoun. Image-based 3d object reconstruction: State-of-the-art and trends in the deep learning era. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 1, 2
- [24] Christian Hane, Shubham Tulsiani, and Jitendra Malik. Hierarchical surface prediction for 3d object reconstruction. In *International Conference on 3D Vision*, pages 412–420, 2017. 2
- [25] Eldar Insafutdinov and Alexey Dosovitskiy. Unsupervised learning of shape and pose with differentiable point clouds. In *Advances in Neural Information Processing Systems*, pages 2807–2817, 2018. 3
- [26] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. SurfacerNet: An end-to-end 3d neural network for multiview stereopsis. In *IEEE International Conference on Computer Vision*, pages 2326–2334, 2017. 2
- [27] Haiyong Jiang, Jianfei Cai, and Jianmin Zheng. Skeleton-aware 3d human shape reconstruction from point clouds. In *IEEE International Conference on Computer Vision*, pages 5431–5441, 2019. 3

- [28] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, 2018. 3
- [29] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. In *Advances in Neural Information Processing Systems*, pages 365–376, 2017. 2
- [30] Michael M. Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Eurographics Symposium on Geometry Processing*, pages 61–70, 2006. 3
- [31] Michael M. Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics*, 32(3):29:1–29:13, 2013. 3
- [32] Nikos Kolotouros, Georgios Pavlakos, Michael J. Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *IEEE International Conference on Computer Vision*, pages 2252–2261, 2019. 3
- [33] Nikos Kolotouros, Georgios Pavlakos, and Kostas Daniilidis. Convolutional mesh regression for single-image human shape reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4501–4510, 2019. 3
- [34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1106–1114, 2012. 4
- [35] Daniel Vlasic Aaron Sarna William T. Freeman Thomas Funkhouser Kyle Genova, Forrester Cole. Learning shape templates with structured implicit functions. Nov 2019. 3
- [36] Lubor Ladicky, Olivier Saurer, SoHyeon Jeong, Fabio Maninchedda, and Marc Pollefeys. From point clouds to mesh using regression. In *IEEE International Conference on Computer Vision*, pages 3913–3922, 2017. 3
- [37] Verica Lazova, Eldar Insafutdinov, and Gerard Pons-Moll. 360-degree textures of people in clothing from a single image. In *International Conference on 3D Vision (3DV)*, sep 2019. 3
- [38] Vincent Leroy, Jean-Sébastien Franco, and Edmond Boyer. Shape reconstruction using volume sweeping and learned photoconsistency. In *European Conference on Computer Vision*, pages 796–811, 2018. 3
- [39] Yiyi Liao, Simon Donné, and Andreas Geiger. Deep marching cubes: Learning explicit surface representations. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2916–2925, 2018. 2, 6
- [40] Shikun Liu, Lee Giles, and Alexander Ororbia. Learning a hierarchical latent-variable model of 3d shapes. In *International Conference on 3D Vision*, pages 542–551. IEEE, 2018. 3
- [41] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM Transactions on Graphics*, 34(6):248:1–248:16, 2015. 3
- [42] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *Computer Graphics and Interactive Techniques*, pages 163–169, 1987. 2, 3, 4, 6
- [43] Lars M. Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019. 2, 3, 4, 6
- [44] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Deep level sets: Implicit surface representations for 3d shape inference. *arXiv preprint arXiv:1901.06802*, 2019. 3
- [45] Richard A. Newcombe, Dieter Fox, and Steven M. Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 343–352, 2015. 3
- [46] Kyle Olszewski, Sergey Tulyakov, Oliver Woodford, Hao Li, and Linjie Luo. Transformable bottleneck networks. *The IEEE International Conference on Computer Vision (ICCV)*, Nov 2019. 2
- [47] Mohamed Omran, Christop Lassner, Gerard Pons-Moll, Peter Gehler, and Bernt Schiele. Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In *International Conference on 3D Vision*, 2018. 3
- [48] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 2, 3, 4
- [49] Chaitanya Patel, Zhouyingcheng Liao, and Gerard Pons-Moll. Tailornet: Predicting clothing in 3d as a function of human pose, shape and garment style. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2020. 3
- [50] Mark Pauly, Markus H. Gross, and Leif Kobbelt. Efficient simplification of point-sampled surfaces. In *IEEE Visualization*, pages 163–170, 2002. 3
- [51] Gerard Pons-Moll, Sergi Pujades, Sonny Hu, and Michael Black. ClothCap: Seamless 4D clothing capture and retargeting. *ACM Transactions on Graphics*, 36(4), 2017. 3
- [52] Sergey Prokudin, Christoph Lassner, and Javier Romero. Efficient learning on point clouds with basis point sets. In *IEEE International Conference on Computer Vision Workshops*, 2019. 3
- [53] Albert Pumarola, Jordi Sanchez, Gary P. T. Choi, Alberto Sanfeliu, and Francesc Moreno-Noguer. 3dpeople: Modeling the geometry of dressed humans. In *IEEE International Conference on Computer Vision*, pages 2242–2251, 2019. 3
- [54] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 77–85, 2017. 2, 3
- [55] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 77–85, 2017. 2, 3
- [56] Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J. Black. Generating 3d faces using convolutional

- mesh autoencoders. In *European Conference on Computer Vision*, pages 725–741, 2018. 3
- [57] Danilo Jimenez Rezende, S. M. Ali Eslami, Shakir Mohamed, Peter W. Battaglia, Max Jaderberg, and Nicolas Heess. Unsupervised learning of 3d structure from images. In *Advances in Neural Information Processing Systems*, pages 4997–5005, 2016. 2
- [58] Gernot Riegler, Ali Osman Ulusoy, Horst Bischof, and Andreas Geiger. Octnetfusion: Learning depth fusion from data. In *International Conference on 3D Vision*, pages 57–66, 2017. 3
- [59] 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 *IEEE International Conference on Computer Vision*, pages 2304–2314, 2019. 4
- [60] Dong Wook Shu, Sung Woo Park, and Junseok Kwon. 3d point cloud generative adversarial network based on tree structured graph convolutions. In *IEEE International Conference on Computer Vision*, pages 3859–3868, 2019. 3
- [61] Miroslava Slavcheva, Maximilian Baust, Daniel Cremers, and Slobodan Ilic. Killingfusion: Non-rigid 3d reconstruction without correspondences. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 3
- [62] David Smith, Matthew Loper, Xiaochen Hu, Paris Mavroidis, and Javier Romero. FACSIMILE: fast and accurate scans from an image in less than a second. In *IEEE International Conference on Computer Vision*, 2019. 3
- [63] Edward Smith, Scott Fujimoto, and David Meger. Multi-view silhouette and depth decomposition for high resolution 3d object representation. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 6479–6489. Curran Associates, Inc., 2018. 2
- [64] David Stutz and Andreas Geiger. Learning 3d shape completion from laser scan data with weak supervision. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1955–1964, 2018. 3
- [65] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *IEEE International Conference on Computer Vision*, pages 2107–2115, 2017. 2
- [66] Maxim Tatarchenko, Stephan R. Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3405–3414, 2019. 2
- [67] 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 *IEEE International Conference on Computer Vision*, pages 4541–4550, 2019. 3
- [68] Shubham Tulsiani, Tinghui Zhou, Alexei A. Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 209–217, 2017. 2
- [69] Hsiao-Yu Tung, Hsiao-Wei Tung, Ersin Yumer, and Katerina Fragkiadaki. Self-supervised learning of motion capture. In *Advances in Neural Information Processing Systems*, pages 5236–5246, 2017. 3
- [70] Gül Varol, Duygu Ceylan, Bryan Russell, Jimei Yang, Ersin Yumer, Ivan Laptev, and Cordelia Schmid. Bodynet: Volumetric inference of 3d human body shapes. In *European Conference on Computer Vision*, 2018. 3
- [71] Hao Wang, Nadav Schor, Ruizhen Hu, Haibin Huang, Daniel Cohen-Or, and Hui Huang. Global-to-local generative model for 3d shapes. *ACM Transactions on Graphics*, 37(6):214:1–214:10, 2018. 2
- [72] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single RGB images. In *European Conference on Computer Vision*, pages 55–71, 2018. 2, 3
- [73] Jiajun Wu, Yifan Wang, Tianfan Xue, Xingyuan Sun, Bill Freeman, and Josh Tenenbaum. Marrnet: 3d shape reconstruction via 2.5d sketches. In *Advances in Neural Information Processing Systems*, pages 540–550, 2017. 2
- [74] 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, 3
- [75] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1912–1920, 2015. 2
- [76] 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. 3, 6
- [77] 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 *IEEE International Conference on Computer Vision*, pages 4541–4550, 2019. 3
- [78] Tao Yu, Zerong Zheng, Kaiwen Guo, Jianhui Zhao, Qionghai Dai, Hao Li, Gerard Pons-Moll, and Yebin Liu. Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 7287–7296, 2018. 3
- [79] Andrei Zanfir, Elisabeta Marinoiu, and Cristian Sminchisescu. Monocular 3d pose and shape estimation of multiple people in natural scenes—the importance of multiple scene constraints. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2148–2157, 2018. 3
- [80] Chao Zhang, Sergi Pujades, Michael Black, and Gerard Pons-Moll. Detailed, accurate, human shape estimation from clothed 3D scan sequences. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 8
- [81] Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Josh Tenenbaum, Bill Freeman, and Jiajun Wu. Learning to re-

construct shapes from unseen classes. In *Advances in Neural Information Processing Systems*, pages 2263–2274, 2018. 2

- [82] Zerong Zheng, Tao Yu, Yixuan Wei, Qionghai Dai, and Yebin Liu. Deephuman: 3d human reconstruction from a single image. In *IEEE International Conference on Computer Vision*, pages 7739–7749, 2019. 3