

3D Hand Shape and Pose Estimation from a Single RGB Image

Liuhao Ge^{1*}, Zhou Ren², Yuncheng Li³, Zehao Xue³, Yingying Wang³, Jianfei Cai¹, Junsong Yuan⁴

¹Nanyang Technological University ²Wormpex AI Research

³Snap Inc. ⁴State University of New York at Buffalo

ge0001ao@e.ntu.edu.sg, zhou.ren@bianlifeng.com, yuncheng.li@snap.com,
zehao.xue@snap.com, ywang@snap.com, asjfcai@ntu.edu.sg, jsyuan@buffalo.edu

Abstract

This work addresses a novel and challenging problem of estimating the full 3D hand shape and pose from a single RGB image. Most current methods in 3D hand analysis from monocular RGB images only focus on estimating the 3D locations of hand keypoints, which cannot fully express the 3D shape of hand. In contrast, we propose a Graph Convolutional Neural Network (Graph CNN) based method to reconstruct a full 3D mesh of hand surface that contains richer information of both 3D hand shape and pose. To train networks with full supervision, we create a large-scale synthetic dataset containing both ground truth 3D meshes and 3D poses. When fine-tuning the networks on real-world datasets without 3D ground truth, we propose a weakly-supervised approach by leveraging the depth map as a weak supervision in training. Through extensive evaluations on our proposed new datasets and two public datasets, we show that our proposed method can produce accurate and reasonable 3D hand mesh, and can achieve superior 3D hand pose estimation accuracy when compared with state-of-the-art methods.

1. Introduction

Vision-based 3D hand analysis is a very important topic because it has many applications in virtual reality (VR) and augmented reality (AR). However, despite years of studies [40, 57, 58, 47, 45, 13, 27], it remains an open problem due to the diversity and complexity of hand shape, pose, gesture, occlusion, *etc.* In the past decade, we have witnessed a rapid advance in 3D hand pose estimation from depth images [35, 52, 12, 15, 14, 61, 11, 16]. Considering RGB cameras are more widely available than depth cameras, some recent works start looking into 3D hand analysis from monocular RGB images, and mainly focus on estimating sparse 3D hand joint locations but ignore dense 3D hand

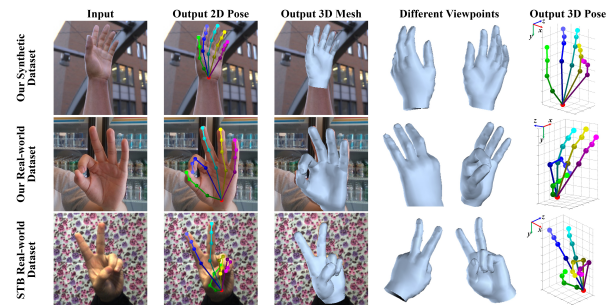


Figure 1: Our proposed method is able to not only estimate 2D/3D hand joint locations, but also recover a full 3D mesh of hand surface from a single RGB image. We show our estimation results on our proposed synthetic and real-world datasets as well as the STB real-world dataset [62].

shape [63, 44, 32, 5, 20, 36, 38]. However, many immersive VR and AR applications often require accurate estimation of both 3D hand pose and 3D hand shape.

This motivates us to bring out a more challenging task: *how to jointly estimate not only the 3D hand joint locations, but also the full 3D mesh of hand surface from a single RGB image?* In this work, we develop a sound solution to this task, as illustrated in Fig. 1.

The task of single-view 3D hand shape estimation has been studied previously, but mostly in controlled settings, where a depth sensor is available. The basic idea is to fit a generative 3D hand model to the input depth image with iterative optimization [49, 30, 24, 21, 51, 41]. In contrast, here we consider to estimate 3D hand shape from a monocular RGB image, which has not been extensively studied yet. The absence of explicit depth cues in RGB images makes this task difficult to be solved by iterative optimization approaches. In this work, we apply deep neural networks that are trained in an end-to-end manner to recover 3D hand mesh directly from a single RGB image. Specifically, we predefine the topology of a triangle mesh representing the hand surface, and aim at estimating the 3D coordinates of all the vertices in the mesh using deep neural networks. To achieve this goal, there are several challenges.

*This work was done when Liuhao Ge was a research intern at Snap Inc.

The first challenge is the high dimensionality of the output space for 3D hand mesh generation. Compared with estimating sparse 3D joint locations of the hand skeleton (*e.g.*, 21 joints), it is much more difficult to estimate 3D coordinates of dense mesh vertices (*e.g.*, 1280 vertices) using conventional CNNs. One straightforward solution is to follow the common approach used in human body shape estimation [53, 48, 37, 22], namely to regress low-dimensional parameters of a predefined deformable hand model, *e.g.*, MANO [42].

In this paper, we argue that the output 3D hand mesh vertices in essence are graph-structured data, since a 3D mesh can be easily represented as a graph. To output such graph-structured data and better exploit the topological relationship among mesh vertices in the graph, motivated by recent works on Graph CNNs [8, 39, 56], we propose a novel Graph CNN-based approach. Specifically, we adopt graph convolutions [8] hierarchically with upsampling and nonlinear activations to generate 3D hand mesh vertices in a graph from image features which are extracted by backbone networks. With such an end-to-end trainable framework, our Graph CNN-based method can better represent the highly variable 3D hand shapes, and can better express the local details of 3D hand shapes.

Besides the computational model, an additional challenge is the lack of ground truth 3D hand mesh training data for real-world images. Manually annotating the ground truth 3D hand meshes on real-world RGB images is extremely laborious and time-consuming. We thus choose to create a large-scale synthetic dataset containing the ground truth of both 3D hand mesh and 3D hand pose for training. However, models trained on the synthetic dataset usually produce unsatisfactory estimation results on real-world datasets due to the domain gap between them. To address this issue, inspired by [5, 37], we propose a novel weakly-supervised method by leveraging depth map as a weak supervision for 3D mesh generation, since depth map can be easily captured by an RGB-D camera when collecting real-world training data. More specifically, when fine-tuning on real-world datasets, we render the generated 3D hand mesh to a depth map on the image plane and minimize the depth map loss against the reference depth map, as shown in Fig. 3. Note that, during testing, we only need an RGB image as input to estimate full 3D hand shape and pose.

To the best of our knowledge, we are the first to handle the problem of estimating not only 3D hand pose but also full 3D hand shape from a single RGB image. Our main contributions are summarized as follows:

- We propose a novel end-to-end trainable hand mesh generation approach based on Graph CNN [8]. Experiments show that our method can well represent hand shape variations and capture local details. Furthermore, we observe that by estimating full 3D hand mesh, our method boost

the accuracy performance of 3D hand pose estimation, as validated in Sec. 5.4.

- We propose a weakly-supervised training pipeline on real-world dataset, by rendering the generated 3D mesh to a depth map on the image plane and leveraging the reference depth map as a weak supervision, without requiring any annotations of 3D hand mesh or 3D hand pose for real-world images.
- We introduce the first large-scale synthetic RGB-based 3D hand shape and pose dataset as well as a small-scale real-world dataset, which contain the annotation of both 3D hand joint locations and the full 3D meshes of hand surface. We will share our datasets publicly upon the acceptance of this work.

We conduct comprehensive experiments on our proposed synthetic and real-world datasets as well as two public datasets [62, 63]. Experimental results show that our proposed method can produce accurate and reasonable 3D hand mesh with real-time speed on GPU, and can achieve superior accuracy performance on 3D hand pose estimation when compared with state-of-the-art methods.

2. Related Work

3D hand shape and pose estimation from depth images:

Most previous methods estimate 3D hand shape and pose from depth images by fitting a deformable hand model to the input depth map with iterative optimization [49, 30, 24, 21, 51, 41]. A recent method [31] was proposed to estimate pose and shape parameters from the depth image using CNNs, and recover 3D hand meshes using LBS. The CNNs are trained in an end-to-end manner with mesh and pose losses. However, the quality of their recovered hand meshes is restricted by their simple LBS model.

3D hand pose estimation from RGB images:

Pioneering works [58, 7] estimate hand pose from RGB image sequences. Gorce et al. [7] proposed estimating 3D hand pose, the hand texture and the illuminant dynamically through minimization of an objective function. Sridhar et al. [46] adopted multi-view RGB images and depth data to estimate the 3D hand pose by combining a discriminative method with local optimization. With the advance of deep learning and the wide applications of monocular RGB cameras, many recent works estimate 3D hand pose from a single RGB image using deep neural networks [63, 44, 32, 5, 20, 38]. However, few works focus on 3D hand shape estimation from RGB images. Panteleris *et al.* [36] proposed to fit a 3D hand model to the estimated 2D joint locations. But the hand model is controlled by 27 hand pose parameters, thus it cannot well adapt to various hand shapes. In addition, this method is not an end-to-end framework for generating 3D hand mesh.

3D human body shape and pose estimation from a single RGB image:

Most recent methods rely on SMPL, a body shape and pose model [29]. Some methods fit the SMPL model to the detected 2D keypoints [3, 25]. Some methods regress SMPL parameters using CNNs with supervisions of silhouette and/or 2D keypoints [48, 37, 22]. A more recent method [54] predicts a volumetric representation of human body. Different from these methods, we propose to estimate 3D mesh vertices using Graph CNNs in order to learn nonlinear hand shape variations and better utilize the relationship among vertices in the mesh topology. In addition, instead of using 2D silhouette or 2D keypoints to weakly supervise the network training, we propose to leverage the depth map as a weak 3D supervision when training on real-world datasets without 3D mesh or 3D pose annotations.

3. 3D Hand Shape and Pose Dataset Creation

Manually annotating the ground truth of 3D hand meshes and 3D hand joint locations for real-world RGB images is extremely laborious and time-consuming. To overcome the difficulties in real-world data annotation, some works [43, 63, 33] have adopted synthetically generated hand RGB images for training. However, existing hand RGB image datasets [43, 62, 63, 33] only provide the annotations of 2D/3D hand joint locations, and they do not contain any 3D hand shape annotations. Thus, these datasets are not suitable for the training of the 3D hand shape estimation task.

In this work, we create a large-scale synthetic hand shape and pose dataset that provides the annotations of both 3D hand joint locations and full 3D hand meshes. In particular, we use Maya [2] to create a 3D hand model and rig it with joints, and then apply photorealistic textures on it as well as natural lighting using High-Dynamic-Range (HDR) images. We model hand variations by creating blend shapes with different shapes and ratios, then applying random weights on the blend shapes. To fully explore the pose space, we create hand poses from 500 common hand gestures and 1000 unique camera viewpoints. To simulate real-world diversity, we use 30 lightings and five skin colors. We render the hand using global illumination with off-the-shelf Arnold renderer [1]. The rendering tasks are distributed onto a cloud render farm for maximum efficiency. In total, our synthetic dataset contains 375,000 hand RGB images with large variations. We use 315,000 images for training and 60,000 images for validation. During training, we randomly sample and crop background images from COCO [28], LSUN [60], and Flickr [10] datasets, and blend them with the rendered hand images, as shown in Fig. 2.

In addition, to quantitatively evaluate the performance of hand mesh estimation on real-world image, we create a real-world dataset containing 583 hand RGB images with the annotations of 3D hand mesh and 3D hand joint locations. To

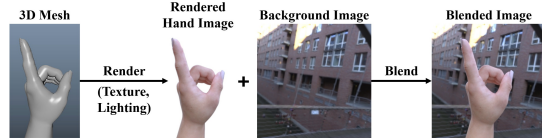


Figure 2: Illustration of our synthetic hand shape and pose dataset creation as well as background image augmentation during training.

facilitate the 3D annotation, we capture the corresponding depth images using an Intel RealSense RGB-D camera [19] and manually adjust the 3D hand model in Maya with the reference of both RGB images and depth points. In this work, this real-world dataset is only used for evaluation.

4. Methodology

4.1. Overview

We propose to generate a full 3D mesh of the hand surface and the 3D hand joint locations directly from a single monocular RGB image, as illustrated in Fig. 3. Specifically, the input is a single RGB image centered on a hand, which is passed through a two-stacked hourglass network [34] to infer 2D heat-maps. The estimated 2D heat-maps, combined with the image feature maps, are encoded as a latent feature vector by using a residual network [18] that contains eight residual layers and four max pooling layers. The encoded latent feature vector is then input to a Graph CNN [8] to infer the 3D coordinates of N vertices $\mathcal{V} = \{v_i\}_{i=1}^N$ in the 3D hand mesh. The 3D hand joint locations $\Phi = \{\phi_j\}_{j=1}^J$ are linearly regressed from the reconstructed 3D hand mesh vertices by using a simplified linear Graph CNN.

In this work, we first train the network models on a synthetic dataset and then fine-tune them on real-world datasets. On the synthetic dataset that contains the ground truth of 3D hand meshes and 3D hand joint locations, we train the networks end-to-end in a fully-supervised manner by using 2D heat-map loss, 3D mesh loss, and 3D pose loss. More details will be presented in Section 4.3. On the real-world dataset, the networks can be fine-tuned in a weakly-supervised manner without requiring the ground truth of 3D hand meshes or 3D hand joint locations. To achieve this target, we leverage the reference depth map available in training, which can be easily captured from a depth camera, as a weak supervision during the fine-tuning, and employ a differentiable renderer to render the generated 3D mesh to a depth map from the camera viewpoint. To guarantee the mesh quality, we generate the pseudo-ground truth mesh from the pretrained model as an additional supervision. More details will be presented in Section 4.4.

4.2. Graph CNNs for Mesh and Pose Estimation

Graph CNNs have been successfully applied in modeling graph structured data [56, 59, 55]. As 3D hand mesh is of

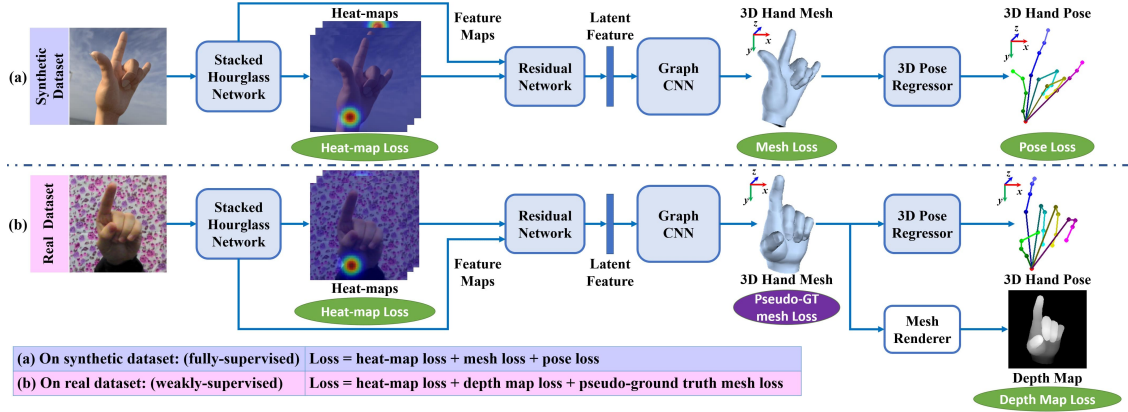


Figure 3: Overview of our method for 3D hand shape and pose estimation from a single RGB image. Our network model is first trained on a synthetic dataset in a fully supervised manner with heat-map loss, 3D mesh loss, and 3D pose loss, as shown in (a); and then fine-tuned on a real-world dataset without 3D mesh or 3D pose ground truth in a weakly-supervised manner by innovatively introducing a pseudo-ground truth mesh loss and a depth map loss, as shown in (b). For both (a) and (b), the input RGB image is first passed through a two-stacked hourglass network [34] for extracting feature maps and 2D heat-maps, which are then combined and encoded as a latent feature vector by a residual network [18]. The latent feature is fed into a Graph CNN [8] to infer the 3D coordinates of mesh vertices. Finally, the 3D hand pose is linearly regressed from the 3D hand mesh. During training on the real-world dataset, as shown in (b), the generated 3D hand mesh is rendered to a depth map to compute the depth map loss against the reference depth map. Note that this step is not involved in testing.

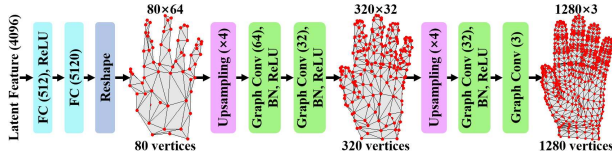


Figure 4: Architecture of the Graph CNN for mesh generation. The input is a latent feature vector extracted from the input RGB image. Passing through two fully-connected (FC) layers, the feature vector is transformed into 80 vertices with 64-dim features in a coarse graph. The features are upsampled and allocated to a finer graph. With two upsampling layers and four graph convolutional layers, the network outputs 3D coordinates of the 1280 mesh vertices. The numbers in parentheses of FC layers and graph convolutions represent the dimensions of output features.

graph structure by nature, in this work we adopt the Chebyshev Spectral Graph CNN [8] to generate 3D coordinates of vertices in the hand mesh and estimate 3D hand pose from the generated mesh.

A 3D mesh can be represented by an undirected graph $\mathcal{M} = (\mathcal{V}, \mathcal{E}, W)$, where $\mathcal{V} = \{\mathbf{v}_i\}_{i=1}^N$ is a set of N vertices in the mesh, $\mathcal{E} = \{\mathbf{e}_i\}_{i=1}^E$ is a set of E edges in the mesh, $W = (w_{ij})_{N \times N}$ is the adjacency matrix, where $w_{ij} = 0$ if $(i, j) \notin \mathcal{E}$, and $w_{ij} = 1$ if $(i, j) \in \mathcal{E}$. The normalized graph Laplacian [6] is computed as $L = I_N - D^{-1/2}WD^{-1/2}$, where $D = \text{diag}(\sum_j w_{ij})$ is the diagonal degree matrix, I_N is the identity matrix. Here, we assume that the topology of the triangular mesh is fixed and is predefined by the hand mesh model, *i.e.*, the adjacency matrix W and the graph Laplacian L of the graph \mathcal{M} are fixed during training and

testing.

Given a signal $\mathbf{f} = (f_1, \dots, f_N)^T \in \mathbb{R}^{N \times F}$ on the vertices of graph \mathcal{M} , it represents F -dim features of N vertices in the 3D mesh. In Chebyshev Spectral Graph CNN [8], the graph convolutional operation on a graph signal $\mathbf{f}_{\text{in}} \in \mathbb{R}^{N \times F_{\text{in}}}$ is defined as

$$\mathbf{f}_{\text{out}} = \sum_{k=0}^{K-1} T_k(\tilde{L}) \cdot \mathbf{f}_{\text{in}} \cdot \theta_k, \quad (1)$$

where $T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x)$ is the Chebyshev polynomial of degree k , $T_0 = 1$, $T_1 = x$; $\tilde{L} \in \mathbb{R}^{N \times N}$ is the rescaled Laplacian, $\tilde{L} = 2L/\lambda_{\max} - I_N$, λ_{\max} is the maximum eigenvalue of L ; $\theta_k \in \mathbb{R}^{F_{\text{in}} \times F_{\text{out}}}$ are the trainable parameters in the graph convolutional layer; $\mathbf{f}_{\text{out}} \in \mathbb{R}^{N \times F_{\text{out}}}$ is the output graph signal. This operation is K -localized since Eq. 1 is a K -order polynomial of the graph Laplacian, and it only affects the K -hop neighbors of each central node. Readers are referred to [8] for more details.

In this work, we design a hierarchical architecture for mesh generation by performing graph convolution on graphs from coarse to fine, as shown in Fig. 4. The topologies of coarse graphs are precomputed by graph coarsening, as shown in Fig. 5 (a), and are fixed during training and testing. Following Defferrard *et al.* [8], we use the Graclus multilevel clustering algorithm [9] to coarsen the graph, and create a tree structure to store correspondences of vertices in graphs at adjacent coarsening levels. During the forward propagation, we upsample features of vertices in the coarse graph to corresponding children vertices in the fine graph, as shown in Fig. 5 (b). Then, we perform the graph convolution to update features in the graph. All the

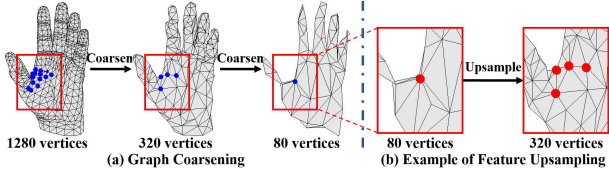


Figure 5: (a) Given our predefined mesh topology, we first perform graph coarsening [8] to cluster meaningful neighborhoods on graphs and create a tree structure to store correspondences of vertices in graphs at adjacent coarsening levels. (b) During the forward propagation, we perform feature upsampling. The feature of a vertex in the coarse graph is allocated to its children vertices in the finer graph.

graph convolutional filters have the same support of $K = 3$. To make the network output irrelevant to the camera intrinsic parameters, we design the network to output UV coordinates on input image and depth of vertices in the mesh, which can be converted to 3D coordinates in the camera coordinate system using the camera intrinsic matrix. Similar to [63, 5, 44], we estimate scale-invariant and root-relative depth of mesh vertices.

Considering that 3D joint locations can be estimated directly from the 3D mesh vertices using a linear regressor [29, 42], we adopt a simplified Graph CNN [8] with two pooling layers and without nonlinear activation to linearly regress the scale-invariant and root-relative 3D hand joint locations from 3D coordinates of hand mesh vertices.

4.3. Fully-supervised Training on Synthetic Dataset

We first train the networks on our synthetic hand shape and pose dataset in a fully-supervised manner. As shown in Fig. 3 (a), the networks are supervised by heat-map loss $\mathcal{L}_{\mathcal{H}}$, mesh loss $\mathcal{L}_{\mathcal{M}}$, and 3D pose loss $\mathcal{L}_{\mathcal{J}}$.

Heat-map Loss. $\mathcal{L}_{\mathcal{H}} = \sum_{j=1}^J \left\| \mathcal{H}_j - \hat{\mathcal{H}}_j \right\|_2^2$, where \mathcal{H}_j and $\hat{\mathcal{H}}_j$ are the ground truth and estimated heat-maps, respectively. We set the heat-map resolution as 64×64 px. The ground truth heat-map is defined as a 2D Gaussian with a standard deviation of 4 px centered on the ground truth 2D joint location.

Mesh Loss. Similar to [56], $\mathcal{L}_{\mathcal{M}} = \lambda_v \mathcal{L}_v + \lambda_n \mathcal{L}_n + \lambda_e \mathcal{L}_e + \lambda_l \mathcal{L}_l$ is composed of vertex loss \mathcal{L}_v , normal loss \mathcal{L}_n , edge loss \mathcal{L}_e , and Laplacian loss \mathcal{L}_l . The vertex loss \mathcal{L}_v is to constrain 2D and 3D locations of mesh vertices:

$$\mathcal{L}_v = \sum_{i=1}^N \left\| \mathbf{v}_i^{3D} - \hat{\mathbf{v}}_i^{3D} \right\|_2^2 + \left\| \mathbf{v}_i^{2D} - \hat{\mathbf{v}}_i^{2D} \right\|_2^2, \quad (2)$$

where \mathbf{v}_i and $\hat{\mathbf{v}}_i$ denote the ground truth and estimated 2D/3D locations of the mesh vertices, respectively. The normal loss \mathcal{L}_n is to enforce surface normal consistency:

$$\mathcal{L}_n = \sum_t \sum_{(i,j) \in t} \left\| \langle \hat{\mathbf{v}}_i^{3D} - \hat{\mathbf{v}}_j^{3D}, \mathbf{n}_t \rangle \right\|_2^2, \quad (3)$$

where t is the index of triangle faces in the mesh; (i, j) are the indices of vertices that compose one edge of triangle t ;

and \mathbf{n}_t is the ground truth normal vector of triangle face t , which is computed from ground truth vertices. The edge loss \mathcal{L}_e is introduced to enforce edge length consistency:

$$\mathcal{L}_e = \sum_{i=1}^E \left(\left\| \mathbf{e}_i \right\|_2^2 - \left\| \hat{\mathbf{e}}_i \right\|_2^2 \right)^2, \quad (4)$$

where \mathbf{e}_i and $\hat{\mathbf{e}}_i$ denote the ground truth and estimated edge vectors, respectively. The Laplacian loss \mathcal{L}_l is introduced to preserve the local surface smoothness of mesh:

$$\mathcal{L}_l = \sum_{i=1}^N \left\| \delta_i - \sum_{\mathbf{v}_k \in \mathcal{N}(\mathbf{v}_i)} \delta_k / B_i \right\|_2^2, \quad (5)$$

where $\delta_i = \mathbf{v}_i^{3D} - \hat{\mathbf{v}}_i^{3D}$ is the offset from the estimation to the ground truth, $\mathcal{N}(\mathbf{v}_i)$ is the set of neighboring vertices of \mathbf{v}_i , and B_i is the number of vertices in the set $\mathcal{N}(\mathbf{v}_i)$. This loss function prevents the neighboring vertices from having opposite offsets, thus making the estimated 3D hand surface mesh smoother. For the hyperparameters, we set $\lambda_v = 1$, $\lambda_n = 1$, $\lambda_e = 1$, $\lambda_l = 50$ in our implementation.

3D Pose Loss. $\mathcal{L}_{\mathcal{J}} = \sum_{j=1}^J \left\| \phi_j^{3D} - \hat{\phi}_j^{3D} \right\|_2^2$, where ϕ_j^{3D} and $\hat{\phi}_j^{3D}$ are the ground truth and estimated 3D joint locations, respectively.

In our implementation, we first train the stacked hourglass network and the 3D pose regressor separately with the heat-map loss and the 3D pose loss, respectively. Then, we train the stacked hourglass network, the residual network and the Graph CNN for mesh generation with the combined loss \mathcal{L}_{fully} :

$$\mathcal{L}_{fully} = \lambda_{\mathcal{H}} \mathcal{L}_{\mathcal{H}} + \lambda_{\mathcal{M}} \mathcal{L}_{\mathcal{M}} + \lambda_{\mathcal{J}} \mathcal{L}_{\mathcal{J}}, \quad (6)$$

where $\lambda_{\mathcal{H}} = 0.5$, $\lambda_{\mathcal{M}} = 1$, $\lambda_{\mathcal{J}} = 1$.

4.4. Weakly-supervised Fine-tuning

On the real-world dataset, *i.e.*, the Stereo Hand Pose Tracking Benchmark [62], there is no ground truth of 3D hand mesh. Thus, we fine-tune the networks in a weakly-supervised manner. Moreover, our model also supports the fine-tuning without the ground truth of 3D joint locations, which can further removes the burden of annotating 3D joint locations on training data and make it more applicable for large-scale real-world dataset.

Depth Map Loss. As shown in Fig. 3 (b), we leverage the reference depth map, which can be easily captured by a depth camera, as a weak supervision, and employ a differentiable renderer, similar to [23], to render the estimated 3D hand mesh to a depth map from the camera viewpoint. We use smooth L1 loss [17] for the depth map loss:

$$\mathcal{L}_{\mathcal{D}} = \text{smooth}_{L1} \left(\mathcal{D}, \hat{\mathcal{D}} \right), \quad \hat{\mathcal{D}} = \mathcal{R} \left(\hat{\mathcal{M}} \right), \quad (7)$$

where \mathcal{D} and $\hat{\mathcal{D}}$ denote the ground truth and rendered depth maps, respectively; $\mathcal{R}(\cdot)$ is the depth rendering function;



Figure 6: Impact of the pseudo-ground truth mesh supervision. Without the supervision of pseudo-ground truth mesh, the network produces very rough meshes with incorrect shape and noisy surface.

$\hat{\mathcal{M}}$ is the estimated 3D hand mesh. We set the resolution of a depth map as 32×32 px.

Pseudo-Ground Truth Mesh Loss. Training with only the depth map loss could lead to a degenerated solution, as shown in Fig. 6 (right), since the depth map loss only constrains the visible surface and is sensitive to the noise in the captured depth map. To solve this issue, inspired by [26], we create the pseudo-ground truth mesh $\tilde{\mathcal{M}}$ by testing on the real-world training data using the pretrained models and the ground truth heat-maps. The pseudo-ground truth mesh $\tilde{\mathcal{M}}$ usually has reasonable edge length and good surface smoothness, although it suffers from the relative depth error. Based on this observation, we do not apply vertex loss or normal loss, and we only adopt the edge loss \mathcal{L}_e and the Laplacian loss \mathcal{L}_l as the pseudo-ground truth mesh loss $\mathcal{L}_{p\mathcal{M}} = \lambda_e \mathcal{L}_e + \lambda_l \mathcal{L}_l$, where $\lambda_e = 1$, $\lambda_l = 50$, in order to preserve the edge length and surface smoothness of the mesh. As shown in Fig. 6 (middle), with the supervision of the pseudo-ground truth meshes, the network can generate meshes with correct shape and smooth surface.

In our implementation, we first fine-tune the stacked hourglass network with the heat-map loss, and then end-to-end fine-tune all networks with the combined loss \mathcal{L}_{weakly} :

$$\mathcal{L}_{weakly} = \lambda_{\mathcal{H}} \mathcal{L}_{\mathcal{H}} + \lambda_{\mathcal{D}} \mathcal{L}_{\mathcal{D}} + \lambda_{p\mathcal{M}} \mathcal{L}_{p\mathcal{M}}, \quad (8)$$

where $\lambda_{\mathcal{H}} = 0.1$, $\lambda_{\mathcal{D}} = 0.1$, $\lambda_{p\mathcal{M}} = 1$. Note that Eq. 8 is the loss function for fine-tuning on the dataset without 3D pose supervision. When the ground truth of 3D joint locations is provided during training, we add the 3D pose loss $\mathcal{L}_{\mathcal{J}}$ in the loss function and set the weight $\lambda_{\mathcal{J}} = 10$.

5. Experiments

5.1. Datasets, Metrics and Implementation Details

In this work, we evaluate our method on two aspects: 3D hand mesh reconstruction and 3D hand pose estimation.

For 3D hand mesh reconstruction, we evaluate the generated 3D hand meshes on our proposed synthetic and real-world datasets, which are introduced in Section 3, since no other hand RGB image dataset contains the ground truth of 3D hand meshes. We measure the average error in Euclidean space between the corresponding vertices in each generated 3D mesh and its ground truth 3D mesh. This metric is denoted as ‘‘mesh error’’ in the following experiments.

Error (mm)	–Normal	–Edge	–Laplacian	–3D Pose	Full
Mesh error	8.34	9.09	8.63	9.04	7.95
Pose error	8.30	9.06	8.55	9.24	8.03

Table 1: Ablation study by eliminating different loss terms from our fully-supervised training loss in Eq. 6, respectively. We report the average mesh and pose errors evaluated on the validation set of our synthetic dataset.

For 3D hand pose estimation, we evaluate our proposed methods on two publicly available datasets: Stereo Hand Pose Tracking Benchmark (STB) [62] and the Rendered Hand Pose Dataset (RHD) [63]. STB is a real-world dataset containing 18,000 images with the ground truth of 21 3D hand joint locations and corresponding depth images. Following [63, 5, 44], we split the dataset into 15,000 training samples and 3,000 test samples. To make the joint definition consistent with our settings and RHD dataset, following [5], we move the root joint location from palm center to wrist. RHD is a synthetic dataset containing 41,258 training images and 2,728 testing images. This dataset is challenging due to the large variations in viewpoints and the low image resolution. We evaluate the performance of 3D hand pose estimation with three metrics: (i) Pose error: the average error in Euclidean space between the estimated 3D joints and the ground truth joints; (ii) 3D PCK: the percentage of correct keypoints of which the Euclidean error distance is below a threshold; (iii) AUC: the area under the curve on PCK for different error thresholds.

We implement our method within the PyTorch framework. The networks are trained using the RMSprop optimizer [50] with mini-batches of size 32. The learning rate is set as 10^{-3} when pretraining on our synthetic dataset, and is set as 10^{-4} when fine-tuning on RHD [63] and STB [62]. The input image is resized to 256×256 px. Following the same condition used in [63, 5, 44], we assume that the global hand scale and the absolute depth of root joint are provided at test time. The global hand scale is set as the length of the bone between MCP and PIP joints of the middle finger.

5.2. Ablation Study of Loss Terms

We first evaluate the impact of different losses used in the fully-supervised training (Eq. 6) on the performance of mesh reconstruction and pose estimation. We conduct this experiment on our synthetic dataset. As presented in Table 1, the model trained with the full loss achieves the best performance in both mesh reconstruction and pose estimation, which indicates that all the losses have contributions to producing accurate 3D hand mesh as well as 3D hand joint locations.

5.3. Evaluation of 3D Hand Mesh Reconstruction

We demonstrate the advantages of our proposed Graph CNN-based 3D hand mesh reconstruction method by com-

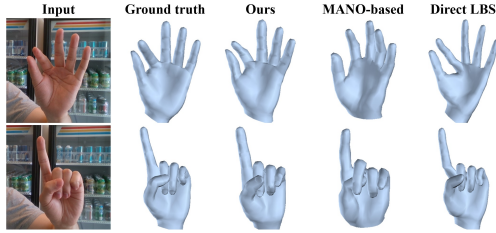


Figure 7: Qualitative comparisons of the meshes generated by our method and other methods. The meshes generated by the MANO-based method usually exhibit inaccurate shape and pose. The meshes generated by the direct Linear Blend Skinning (LBS) method suffer from serious artifacts. Examples are taken from our real-world dataset.

Mesh error (mm)	MANO-based	Direct LBS	Ours
Our synthetic dataset	12.12	10.32	8.01
Our real-world dataset	20.86	13.33	12.72

Table 2: Average mesh errors tested on the validation set of our synthetic dataset and our real-world dataset. We compare our method with two baseline methods. Note that the mesh errors in this table are measured on the aligned mesh defined by MANO [42] for fair comparison.

paring it with two baseline methods: direct Linear Blend Skinning (LBS) method and MANO-based method.

Direct LBS. We train the network to directly regress 3D hand joint locations from the heat-maps and the image features, which is similar to the network architecture proposed in [5]. We generate the 3D hand mesh from only the estimated 3D hand joint locations by applying inverse kinematics and LBS with the predefined mesh model and skinning weights (see the supplementary for details). As shown in Table 2, the average mesh error of direct LBS method is worse than our method on both our synthetic dataset and our real-world dataset, since the LBS model for mesh generation is predefined and cannot be adapt to hands with different shapes. As can be seen in Fig. 7, the hand meshes generated by direct LBS method have unrealistic deformation at joints and suffer from serious inherent artifacts.

MANO-based Method. We also implement a MANO [42] based method that regresses hand shape and pose parameters from the latent image features using three fully-connected layers. Then, the 3D hand mesh is generated from the estimated shape and pose parameters using MANO hand model [42] (see the supplementary for details). The networks are trained in fully-supervised manner using the same loss functions as Eq. 6 on our synthetic dataset. For fair comparison, we align our hand mesh with the MANO hand mesh, and compute mesh error on the aligned mesh. As shown in Table 2 and Fig. 7, the MANO-based method exhibits inferior performance on mesh reconstruction compared with our method. Note that direct supervising MANO parameters on synthetic dataset may obtain better perfor-

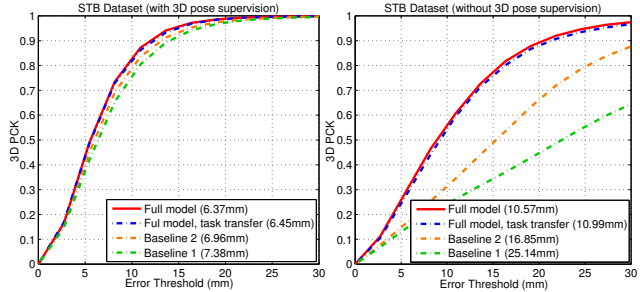


Figure 8: Self-comparisons of 3D hand pose estimation on STB dataset [62]. **Left:** 3D PCK of the model fine-tuned with 3D hand pose supervision. **Right:** 3D PCK of the model fine-tuned without 3D hand pose supervision. The average pose errors are shown in parentheses.

Method	Pipeline	Depth map loss
Baseline 1	im→hm+feat→pose	✗
Baseline 2	im→hm+feat→mesh→pose	✗
Full model	im→hm+feat→mesh→pose	✓

Table 3: Differences between the baseline methods for 3D hand pose estimation and our full model.

mance [4]. But it is infeasible on our synthetic dataset since our dataset does not contain MANO parameters.

5.4. Evaluation of 3D Hand Pose Estimation

We also evaluate our approach on the task of 3D hand pose estimation.

Self-comparisons. We conduct self-comparisons on STB dataset [62] by fine-tuning the networks pretrained on our synthetic dataset in a weakly-supervised manner, as described in Section 4.4. In Table 3, we compare our proposed weakly-supervised method (**Full model**) with two baselines: (i) **Baseline 1:** directly regressing 3D hand joint locations from the heat-maps and the feature maps without using the depth map loss during training; (ii) **Baseline 2:** regressing 3D hand joint locations from the estimated 3D hand mesh without using the depth map loss during training. As presented in Fig. 8, the estimation accuracy of Baseline 2 is superior to that of Baseline 1, which indicates that our proposed 3D hand mesh reconstruction network is beneficial to 3D hand pose estimation. Furthermore, the estimation accuracy of our full model is superior to that of Baseline 2, especially when fine-tuning without 3D hand pose supervision, which validates the effectiveness of introducing the depth map loss as a weak supervision.

In addition, to explore a more efficient way for 3D hand pose estimation without mesh generation, we directly regress the 3D hand joint locations from the latent feature extracted by our full model instead of regressing them from the 3D hand mesh (see the supplementary for details). This task transfer method is denoted as “**Full model, task transfer**” in Fig. 8. Although this method has the same pipeline

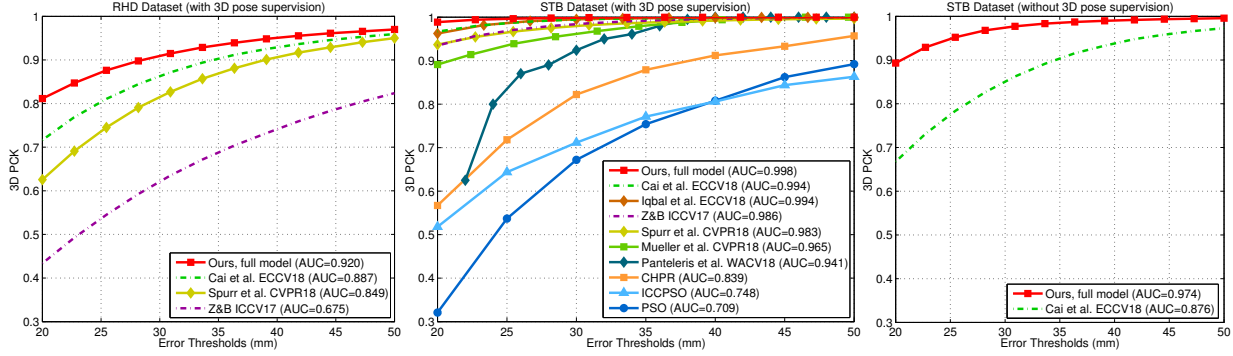


Figure 9: Comparisons with state-of-the-art methods on RHD [63] and STB [62] dataset. **Left:** 3D PCK on RHD dataset [63] with 3D hand pose supervision. **Middle:** 3D PCK on STB dataset [62] with 3D hand pose supervision. **Right:** 3D PCK on STB dataset [62] without 3D hand pose supervision. The AUC values are shown in parentheses.

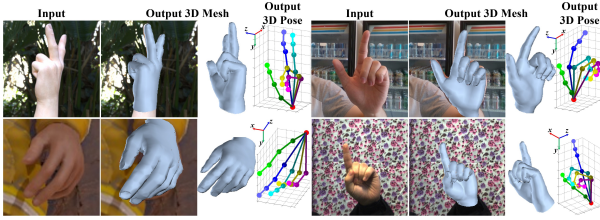


Figure 10: Qualitative results for our synthetic dataset (top left), our real-world dataset (top right), RHD dataset [63] (bottom left), and STB dataset [62] (bottom right).

as that of Baseline 1, the estimation accuracy of this task transfer method is better than that of Baseline 1 and is only a little bit worse than that of our full model, which indicates that the latent feature extracted by our full model is more discriminative and is easier to regress accurate 3D hand pose than the latent feature extracted by Baseline 1.

Comparisons with State-of-the-arts. We compare our method with state-of-the-art 3D hand pose estimation methods on RHD [63] and STB [62] datasets. The PCK curves over different error thresholds are presented in Fig. 9. On RHD dataset, as shown in Fig. 9 (left), our method outperforms the three state-of-the-art methods [63, 44, 5] over all the error thresholds on this dataset. On STB dataset, when the 3D hand pose ground truth is given during training, we compare our methods with seven state-of-the-art methods [62, 63, 36, 44, 32, 5, 20], and our method outperforms these methods over most of the error thresholds, as shown in Fig. 9 (middle). We also experiment with the situation when 3D hand pose ground truth is unknown during training on STB dataset, and compare our method with the weakly-supervised method proposed by Cai *et al.* [5], both of which adopt reference depth maps as a weak supervision. As shown in Fig. 9 (right), our 3D mesh-based method outperforms Cai *et al.* [5] by a large margin.

5.5. Runtime and Qualitative Results

Runtime. We evaluate the runtime of our method on one Nvidia GTX 1080 GPU. The runtime of our full model

outputting both 3D hand mesh and 3D hand pose is 19.9ms on average, including 12.6ms for the stacked hourglass network forward propagation, 4.7ms for the residual network and Graph CNN forward propagation, and 2.6ms for the forward propagation of the pose regressor. Thus, our method can run in real-time on GPU at over 50fps.

Qualitative Results. Some qualitative results of 3D hand mesh reconstruction and 3D hand pose estimation for our synthetic dataset, our real-world dataset, RHD [63], and STB [62] datasets are shown in Fig. 10. More qualitative results are presented in the supplementary.

6. Conclusion

In this paper we have tackled the challenging task of 3D hand shape and pose estimation from a single RGB image. We have developed a Graph CNN-based model to reconstruct a full 3D mesh of hand surface from an input RGB image. To train the model, we have created a large-scale synthetic RGB image dataset with ground truth annotations of both 3D joint locations and 3D hand meshes, on which we train our model in a fully-supervised manner. To fine-tune our model on real-world datasets without 3D ground truth, we render the generated 3D mesh to a depth map and leverage the observed depth map as a weak supervision. Experiments on our proposed new datasets and two public datasets show that our method can recover accurate 3D hand mesh and 3D joint locations in real-time.

In future work, we will use Mocap data to create a larger 3D hand pose and shape dataset. We will also consider the cases of hand-object and hand-hand interactions in order to make the hand pose and shape estimation more robust.

Acknowledgment: This work is in part supported by MoE Tier-2 Grant (2016-T2-2-065) of Singapore. This work is also supported in part by start-up grants from University at Buffalo and a gift grant from Snap Inc.

References

- [1] Autodesk. Arnold renderer. <https://www.arnoldrenderer.com>, 2018.
- [2] Autodesk. Maya. <https://www.autodesk.com.sg/products/maya>, 2018.
- [3] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it simple: Automatic estimation of 3d human pose and shape from a single image. In *ECCV*, 2016.
- [4] Adnane Boukhayma, Rodrigo de Bem, and Philip HS Torr. 3d hand shape and pose from images in the wild. *CVPR*, 2019.
- [5] Yujun Cai, Lihao Ge, Jianfei Cai, and Junsong Yuan. Weakly-supervised 3d hand pose estimation from monocular rgb images. In *ECCV*, 2018.
- [6] Fan RK Chung and Fan Chung Graham. *Spectral graph theory*, volume 92. American Mathematical Society, 1997.
- [7] Martin de La Gorce, David J Fleet, and Nikos Paragios. Model-based 3d hand pose estimation from monocular video. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(9):1793–1805, 2011.
- [8] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *NIPS*, 2016.
- [9] Inderjit S Dhillon, Yuqiang Guan, and Brian Kulis. Weighted graph cuts without eigenvectors a multilevel approach. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(11), 2007.
- [10] Flickr. Flickr. <https://www.flickr.com/>, 2018.
- [11] Lihao Ge, Yujun Cai, Junwu Weng, and Junsong Yuan. Hand pointnet: 3d hand pose estimation using point sets. In *CVPR*, 2018.
- [12] Lihao Ge, Hui Liang, Junsong Yuan, and Daniel Thalmann. Robust 3D hand pose estimation in single depth images: from single-view CNN to multi-view CNNs. In *CVPR*, 2016.
- [13] Lihao Ge, Hui Liang, Junsong Yuan, and Daniel Thalmann. 3D convolutional neural networks for efficient and robust hand pose estimation from single depth images. In *CVPR*, 2017.
- [14] Lihao Ge, Hui Liang, Junsong Yuan, and Daniel Thalmann. Real-time 3D hand pose estimation with 3d convolutional neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [15] Lihao Ge, Hui Liang, Junsong Yuan, and Daniel Thalmann. Robust 3d hand pose estimation from single depth images using multi-view cnns. *IEEE Transactions on Image Processing*, 27(9):4422–4436, 2018.
- [16] Lihao Ge, Zhou Ren, and Junsong Yuan. Point-to-point regression pointnet for 3d hand pose estimation. In *ECCV*, 2018.
- [17] Ross Girshick. Fast r-cnn. In *ICCV*, 2015.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [19] Intel. Intel realsense. <https://realsense.intel.com/>, 2018.
- [20] Umar Iqbal, Pavlo Molchanov, Thomas Breuel, Juergen Gall, and Jan Kautz. Hand pose estimation via latent 2.5D heatmap regression. In *ECCV*, 2018.
- [21] David Joseph Tan, Thomas Cashman, Jonathan Taylor, Andrew Fitzgibbon, Daniel Tarlow, Sameh Khamis, Shahram Izadi, and Jamie Shotton. Fits like a glove: Rapid and reliable hand shape personalization. In *CVPR*, 2016.
- [22] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018.
- [23] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *CVPR*, 2018.
- [24] Sameh Khamis, Jonathan Taylor, Jamie Shotton, Cem Keskin, Shahram Izadi, and Andrew Fitzgibbon. Learning an efficient model of hand shape variation from depth images. In *CVPR*, 2015.
- [25] Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J Black, and Peter V Gehler. Unite the people: Closing the loop between 3d and 2d human representations. In *CVPR*, 2017.
- [26] Zhizhong Li and Derek Hoiem. Learning without forgetting. In *ECCV*, 2017.
- [27] Hui Liang, Junsong Yuan, Jun Lee, Lihao Ge, and Daniel Thalmann. Hough forest with optimized leaves for global hand pose estimation with arbitrary postures. *IEEE Transactions on Cybernetics*, 49(2):527–541, 2019.
- [28] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- [29] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM Transactions on Graphics (TOG)*, 34(6):248, 2015.
- [30] Alexandros Makris and A Argyros. Model-based 3d hand tracking with on-line hand shape adaptation. *BMVC*, 2015.
- [31] Jameel Malik, Ahmed Elhayek, Fabrizio Nunnari, Kiran Varanasi, Kiarash Tamaddon, Alexis Heloir, and Didier Stricker. DeepHps: End-to-end estimation of 3d hand pose and shape by learning from synthetic depth. In *3DV*, 2018.
- [32] Franziska Mueller, Florian Bernard, Oleksandr Sotnychenko, Dushyant Mehta, Srinath Sridhar, Dan Casas, and Christian Theobalt. GANerated hands for real-time 3d hand tracking from monocular RGB. In *CVPR*, 2018.
- [33] Franziska Mueller, Dushyant Mehta, Oleksandr Sotnychenko, Srinath Sridhar, Dan Casas, and Christian Theobalt. Real-time hand tracking under occlusion from an egocentric RGB-D sensor. In *ICCV*, 2017.
- [34] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016.
- [35] Iason Oikonomidis, Nikolaos Kyriazis, and Antonis Argyros. Efficient model-based 3D tracking of hand articulations using Kinect. In *BMVC*, 2011.
- [36] Paschalis Panteleris, Iason Oikonomidis, and Antonis Argyros. Using a single rgb frame for real time 3d hand pose estimation in the wild. In *WACV*, 2018.

- [37] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. *CVPR*, 2018.
- [38] Mahdi Rad, Markus Oberweger, and Vincent Lepetit. Domain transfer for 3d pose estimation from color images without manual annotations. In *ACCV*, 2018.
- [39] Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J Black. Generating 3d faces using convolutional mesh autoencoders. *ECCV*, 2018.
- [40] James M Rehg and Takeo Kanade. Visual tracking of high d-of articulated structures: an application to human hand tracking. In *ECCV*, 1994.
- [41] Edoardo Remelli, Anastasia Tkach, Andrea Tagliasacchi, and Mark Pauly. Low-dimensionality calibration through local anisotropic scaling for robust hand model personalization. In *ICCV*, 2017.
- [42] Javier Romero, Dimitrios Tzionas, and Michael J Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics (TOG)*, 36(6):245, 2017.
- [43] Tomas Simon, Hanbyul Joo, Iain A Matthews, and Yaser Sheikh. Hand keypoint detection in single images using multi-view bootstrapping. In *CVPR*, 2017.
- [44] Adrian Spurr, Jie Song, Seonwook Park, and Otmar Hilliges. Cross-modal deep variational hand pose estimation. In *CVPR*, 2018.
- [45] Srinath Sridhar, Franziska Mueller, Michael Zollhöfer, Dan Casas, Antti Oulasvirta, and Christian Theobalt. Real-time joint tracking of a hand manipulating an object from rgb-d input. In *ECCV*, 2016.
- [46] Srinath Sridhar, Antti Oulasvirta, and Christian Theobalt. Interactive markerless articulated hand motion tracking using rgb and depth data. In *ICCV*, 2013.
- [47] Björn Stenger, Arasanathan Thayananthan, Philip HS Torr, and Roberto Cipolla. Model-based hand tracking using a hierarchical bayesian filter. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(9):1372–1384, 2006.
- [48] Vince Tan, Ignas Budvytis, and Roberto Cipolla. Indirect deep structured learning for 3d human body shape and pose prediction. In *BMVC*, 2017.
- [49] Jonathan Taylor, Richard Stebbing, Varun Ramakrishna, Cem Keskin, Jamie Shotton, Shahram Izadi, Aaron Hertzmann, and Andrew Fitzgibbon. User-specific hand modeling from monocular depth sequences. In *CVPR*, 2014.
- [50] Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop, coursera: Neural networks for machine learning. *University of Toronto, Technical Report*, 2012.
- [51] Anastasia Tkach, Andrea Tagliasacchi, Edoardo Remelli, Mark Pauly, and Andrew Fitzgibbon. Online generative model personalization for hand tracking. *ACM Transactions on Graphics (TOG)*, 36(6):243, 2017.
- [52] Jonathan Tompson, Murphy Stein, Yann Lecun, and Ken Perlin. Real-time continuous pose recovery of human hands using convolutional networks. *ACM Transactions on Graphics (ToG)*, 33(5):169, 2014.
- [53] Hsiao-Yu Tung, Hsiao-Wei Tung, Ersin Yumer, and Katerina Fragkiadaki. Self-supervised learning of motion capture. In *NIPS*, 2017.
- [54] Gül Varol, Duygu Ceylan, Bryan Russell, Jimei Yang, Ersin Yumer, Ivan Laptev, and Cordelia Schmid. Bodynet: Volumetric inference of 3d human body shapes. *ECCV*, 2018.
- [55] Nitika Verma, Edmond Boyer, and Jakob Verbeek. Feastnet: Feature-steered graph convolutions for 3d shape analysis. In *CVPR*, 2018.
- [56] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *ECCV*, 2018.
- [57] Ying Wu and Thomas S Huang. Hand modeling, analysis and recognition. *IEEE Signal Processing Magazine*, 18(3):51–60, 2001.
- [58] Ying Wu, John Lin, and Thomas S Huang. Analyzing and capturing articulated hand motion in image sequences. *IEEE transactions on pattern analysis and machine intelligence*, 27(12):1910–1922, 2005.
- [59] Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. *arXiv preprint arXiv:1801.07455*, 2018.
- [60] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.
- [61] Shanxin Yuan, Guillermo Garcia-Hernando, Björn Stenger, Gyeongsik Moon, Ju Yong Chang, Kyoung Mu Lee, Pavlo Molchanov, Jan Kautz, Sina Honari, Liuhao Ge, et al. Depth-based 3d hand pose estimation: From current achievements to future goals. In *CVPR*, 2018.
- [62] Jiawei Zhang, Jianbo Jiao, Mingliang Chen, Liangqiong Qu, Xiaobin Xu, and Qingxiang Yang. 3D hand pose tracking and estimation using stereo matching. *arXiv preprint arXiv:1610.07214*, 2016.
- [63] Christian Zimmermann and Thomas Brox. Learning to estimate 3d hand pose from single RGB images. In *ICCV*, 2017.