

HandVoxNet: Deep Voxel-Based Network for 3D Hand Shape and Pose Estimation from a Single Depth Map

Jameel Malik^{1,2,3}
Sk Aziz Ali^{1,2}

Ibrahim Abdelaziz^{1,2}
Vladislav Golyanik⁵

Ahmed Elhayek^{2,4}
Christian Theobalt⁵

Soshi Shimada⁵
Didier Stricker^{1,2}

¹TU Kaiserslautern ²DFKI Kaiserslautern ³NUST Pakistan ⁴UPM Saudi Arabia ⁵MPII Saarland

Hybrid system

Problem

Directly regress 3D hand meshes from 2D depth images via 2D convolutional neural networks cause perspective distortions

Solution

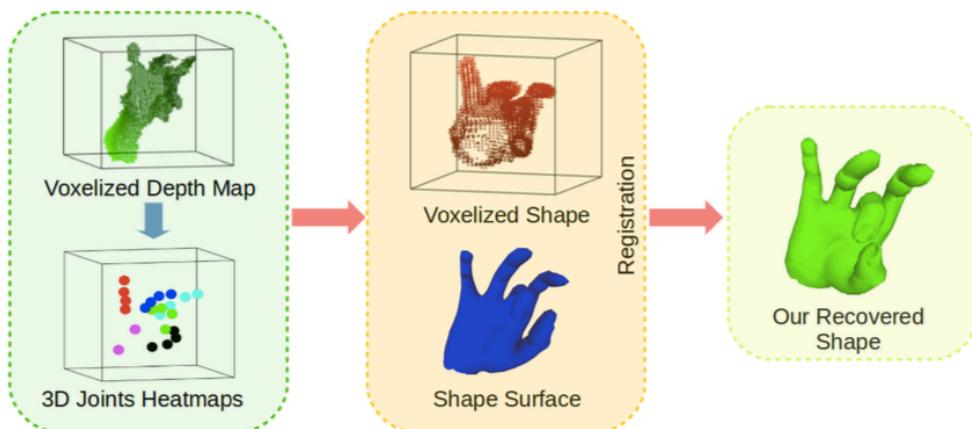
Input 3D voxelized depth map, input depth image is converted into a voxelized grid (*i.e.*, VD) of size $88 \times 88 \times 88$, by using intrinsic camera parameters and a fixed cube size

Using 2 hand shape representations:

3D voxelized grid of the shape

3D hand surface

Combine from registering the hand surface to the voxelized hand shape



Output

N 3D hand joint locations $J \in R^{3 \times N}$ (*i.e.*, 3D pose)

$K = 1193$ 3D vertex locations $V \in R^{3 \times K}$ (*i.e.*, 3D shape).

Solution

Register the estimated shape by V2S-Net to the probabilistic shape representation estimated by FCN (V2V-ShapeNet) using DispVoxNets pipeline

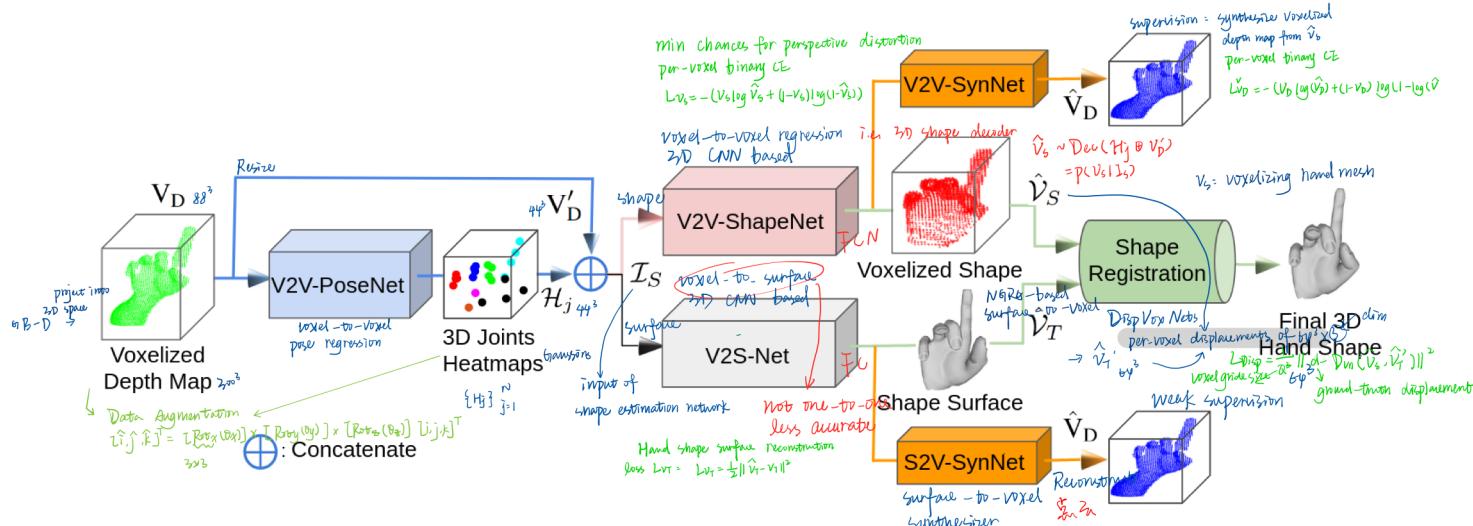


Figure 2: Overview of our approach for 3D hand shape and pose recovery from a 3D voxelized depth map. V2V-PoseNet estimates 3D joints heatmaps (i.e., pose). Hand shape is obtained in two phases. First, V2V-ShapeNet and V2S-Net estimate the voxelized shape and shape surface, respectively. Thereby, V2V-SynNet and S2V-SynNet synthesize the voxelized depth acting as sources of weak-supervision. They are excluded during testing. In the second phase, shape registration accurately fits the shape surface to the voxelized shape.

	Good point	Drawback
V2V-ShapeNet(FC)	Estimate hand shapes while preserving the order and number of points	Local special info loss
V2S-Net(FCN)	Geometry regression	Inconsistent number of points and loses point order.

Data Augmentation in 3D

Originally treat depth maps as 2D data; The representation of the depth map in voxelized form makes it convenient to perform data augmentation in all three dimensions.

Understanding Human Hands in Contact at Internet Scale

Dandan Shan¹, Jiaqi Geng^{*1}, Michelle Shu^{*2}, David F. Fouhey¹

¹University of Michigan, ²Johns Hopkins University

{dandans, jiaqig, fouhey}@umich.edu, msh1@jhu.edu

Detect hand box in RGB that enables mesh reconstruction systems => hand detections

数据采集

Gathering implicit video dataset

identifying an overcomplete set of candidate videos using generic queries and filtering out irrelevant videos.

Annotation: For every hand in each image, we obtained the following annotations: (a) a bounding box around the hand; (b) side: left / right, which is crucial for mesh re- construction; (c) the hand contact state ({no contact, self-

contact, other person contact, in contact with portable object, in contact with a non-portable object}), which provides insights into what the person is doing; and (d) a bounding box around the object the person is contacting *irrespective of name*.

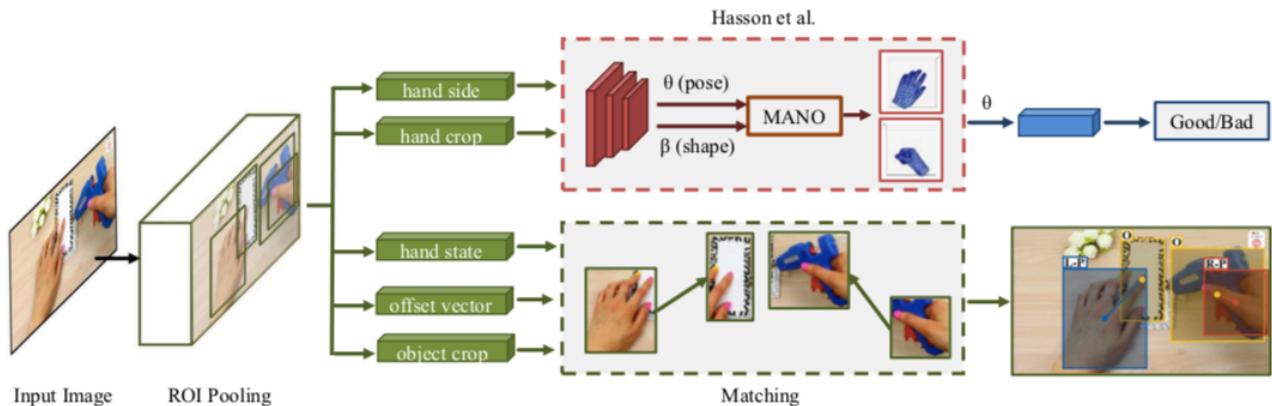


Figure 3. Our system can act as a foundation to understand interacting hands on the Internet. Our system takes a single RGB image and detects hands (irrespective of scale) and for every hand predicts: a box, side, contact state, and a box around the object it is touching. We can then (1) obtain a parse of hand state; and (2) use the hand box and side to feed a reconstruction system like [21]. To help make better use of Internet reconstructions, we introduce a self-supervised system that assesses mesh quality that we train on our data.

Weakly-Supervised Mesh-Convolutional Hand Reconstruction in the Wild

Dominik Kulon^{1, 3} Riza Alp Güler^{1, 3}
 Iasonas Kokkinos³ Michael Bronstein^{1, 2, 4} Stefanos Zafeiriou^{1, 3}

¹Imperial College London ²USI Lugano ³Ariel AI ⁴Twitter

{d.kulon17, r.guler, m.bronstein, s.zafeiriou}@imperial.ac.uk iasonas@arielai.com

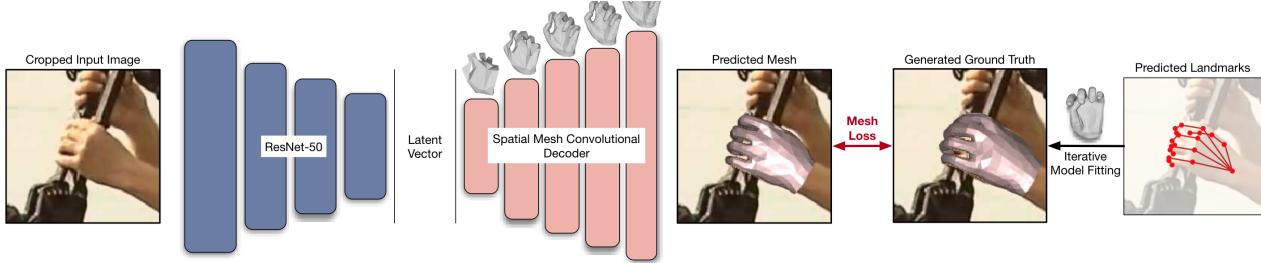


Figure 1: We propose an approach for end-to-end neural network training with mesh supervision that is obtained through an automated data collection method. We process a large collection of YouTube videos and analyze them with 2D hand keypoint detector followed by parametric model fitting (right side). The fitting results are used as a supervisory signal ('mesh loss') for a feed-forward network with a mesh convolutional decoder tasked with recovering a 3D hand mesh at its output (left side).

AUTOMATED DATASET GENERATION FOR 3D RECONSTRUCTION AND POSE ESTIMATION

Detect key points using OpenPose

Lift key points into **3D shapes** by iteratively fitting a **3D deformable model**

3D Shape Representation

$$M(\beta, \theta, \vec{T}_\delta, s; \phi) : \mathbb{R}^{|\beta| \times |\theta| \times |\vec{T}_\delta| \times |s|} \rightarrow \mathbb{R}^{N \times 3} \quad (1)$$

$$\theta_i = P(w)_i = \frac{\sum_{c=1}^C \exp(w^c) P_i^c}{\sum_{c=1}^C \exp(w^c)}.$$

↑ pre-computed
cluster(2) centers

Total #clusters $C = 64$

Represent all constrained angles in terms of w in $R^{K \times C}$

Does not model pairwise dependencies => future work

Requires only a small dataset of angles

Parametric Model Fitting

Regress a matching 3D pose ($K = 16$ joint and $F = 5$ fingertip positions) from MANO mesh #N vertices

$$J(\beta, w, \vec{T}_\delta, s) = \mathcal{J}^T M(\beta, P(w), \vec{T}_\delta, s). \quad (3)$$

$\overset{N \times 3}{\underset{N \times (K+F)}{}}$

Fit the model to 2D annotation by objective:

$$\{\beta^*, w^*, \vec{T}_\delta^*, s^*\} = \arg \min_{\beta, w, \vec{T}_\delta, s} (E_{2D} + E_{bone} + E_{reg}), \quad (4)$$

↑ 2D projection ↑ Joint preservation
 ↓ regularization

used for weak supervision – ground truth

2D landmarks extracted by OpenPose

HAND RECONSTRUCTION NETWORK

Spiral Operator

Only consider $k=2$ here.

$$\begin{aligned} 0\text{-ring}(v) &= \{v\}, \\ (k+1)\text{-ring}(v) &= N(k\text{-ring}(v)) \setminus k\text{-disk}(v), \\ k\text{-disk}(v) &= \bigcup_{i=0 \dots k} i\text{-ring}(v), \end{aligned} \quad (7)$$

$$\mathbf{S}(v) = (v, 1\text{-ring}(v), \dots, k\text{-ring}(v)). \quad (8)$$

$$(f * g)_v = \sum_{\ell=1}^L g_\ell f(S_\ell(v)), \quad (9)$$

1st neighbour
 L kernel
 spiral conv
 spiral length, fixed in each layer

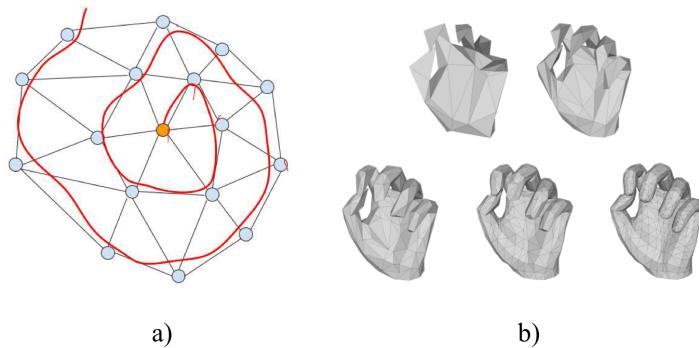


Figure 2: a) Spiral selection of neighbours for a center vertex in clockwise order (taken from [6]). b) The hierarchy of topologies used in our decoder.

Down sampling

#vertices reduced by the factor of 2

Contract vertex pairs based on **quadric error metrics**,

- 1 Project each collapsed node $v_q \in V$ into the closest triangle $v_i, v_j, v_k \in V_d$ in the downsampled mesh V_d obtaining \tilde{v}_p .
- 2 Use the barycentric coordinates w_1, w_2, w_3 such that $\tilde{v}_p = w_1 v_i + w_2 v_j + w_3 v_k$ and $w_1 + w_2 + w_3 = 1$ of the projected vertex to define interpolation weights for the upsampling matrix $Q_u \in R_{m \times n}$,

Up sampling

$V_u = Q_u V_d$ and $m > n$ hold, is formed by setting $Q_u(q, i) = w_i$, $Q_u(q, j) = w_j$, $Q_u(q, k) = w_k$, and $Q_u(q, l) = 0$ for $l \notin \{i, j, k\}$.

The number of vertices at each layer $n \in \{51, 100, 197, 392, 778\}$.

Architecture Encoder-Decoder

Image crop X

=> latent vector $Z = E_{image}(X)$; encoder E_{image} ResNet-50

=> mesh $y = D_{mesh}(Z)$; spiral decoder D_{mesh}



, leaking

ReLU as activation function; Decode a mesh directly from the image encoding and apply spatial convolutions with pooling layers

- 1 Learn invariant shape features in a triangular mesh
- 2 Apply spiral filters to generate hands directly from an image encoding (spatial convolution)

Training

$$\begin{aligned} \mathcal{L} = & \lambda_{vertex} |\hat{\mathcal{Y}} - \mathcal{Y}|_1 \\ & + \lambda_{edge} \sum_{(u,v) \in \mathcal{E}_{mesh}} | ||\hat{\mathcal{Y}}_v - \hat{\mathcal{Y}}_u|| - ||\mathcal{Y}_v - \mathcal{Y}_u|| | \end{aligned} \quad (10)$$

No explicit pose estimation loss; joint coordinates are obtained from eq.3

The network is trained with Adam optimizer

Adapts HOI domain to the single-HPE domain in **image space** rather than feature space via 2 **image generation** methods:

- 1 2D monocular guidance by **GAN** (Generative adversarial network) => align hands
- 2 3D mesh guidance by **MR** (mesh renderer) using estimated 3D meshes and textures => fill in occluded pixels

Thus 1) 3D HPE more accurate 2) HOI input images are translated to segmented and de-occluded hand-only images => more accurate

The pipeline is trained by hand-only images with pose labels and HOI images without pose labels

Simultaneously: 122

- 1 domain adaptation: weak supervision by 2D object segmentation masks and 3D pose (hand-only) labels
- 2 HPE

Problem: lack HOI 3D annotations

Solution: training a new HPE under the HOI scenario by mapping input HOI image to the corresponding hand-only image thus use only datasets:

- 1 Input RGB in hand-only and HOI
 - 2 Skeleton annotations for hand-only
 - 3 2D binary segmentation masks for hand-only and HOI (extracted based on depth maps)
- ⇒ Require restoring occluded regions

$D_{\text{Hand}}^R = \{(x, s_{\text{Hand}}, y)\}$	Real hand-only data (<i>STB</i>)
$D_{\text{Hand}}^S = \{(x, s_{\text{Hand}}, y)\}$	Synthetic hand-only data (<i>SynthHands, RHD</i>)
$D_{\text{HOI}} = \{(x, s_{\text{HOI}})\}$	Real HOI data (<i>CORe50</i>)
$D_{\text{Paired}}^S = \{(x, x_*, s_{\text{Hand}}, s_{\text{HOI}})\}$	Paired synthetic HOI (x) and hand-only (x_*) images (<i>Obman</i>)
$D_{\text{Hand}} = [D_{\text{Hand}}^R, D_{\text{Hand}}^S]$	Hand-only data
$D = [D_{\text{Hand}}, D_{\text{HOI}}^R, D_{\text{Paired}}^S]$	All training data that we use

$X \subset \mathbb{R}^{256 \times 256 \times 3}$	RGB images (x : input; x' : rendered by g^{MR} ; x'' : synthesized by g^{GAN} ; z : final mesh estimate rendered by g^{MR})
$Y \subset \mathbb{R}^{21 \times 3}$	3D skeletal pose space
$F \subset \mathbb{R}^{128 \times 32 \times 32}$	2D feature space
$H \subset \mathbb{R}^{21 \times 32 \times 32}$	2D heatmap space <i>estimated y in image plane</i>
$M \subset \mathbb{R}^{778 \times 1538}$	3D mesh space: 778 vertices \times 1,538 faces
$T \subset \mathbb{R}^{1538 \times 3}$	RGB mesh texture ($3 \times 1,538$ faces)
$g^{\text{FPE}} : X \rightarrow F \times H$	2D feature and pose estimator <i>extracts 2D spatial feature</i>
$g^{\text{HME}} : F \times H \rightarrow M$	Hand mesh estimator <i>maps f & H and generates heatmap h</i>
$g^{\text{Tex}} : F \times H \rightarrow T$	Texture estimator <i>estimates S</i> \rightarrow <i>outputs {f, h, T}</i>
$g^{\text{NR}} : M \times T \rightarrow X$	Neural renderer [25k] <i>maps M to Texture</i>
$g^{\text{Reg}} : M \rightarrow Y$	Hand joint regressor [54] <i>synthesizes x'</i> , x
$g^{\text{MR}} : F \times H \rightarrow X \times Y$	Mesh renderer: $g^{\text{MR}} = [g^{\text{NR}} \circ [g^{\text{HME}}, g^{\text{Tex}}], g^{\text{Reg}} \circ g^{\text{HME}}]$
$g^{\text{GAN}} : F \times H \times F \times H \rightarrow X$	GAN generator <i>refined x' Calculate y from m</i>
$d_1^{\text{GAN}}, d_2^{\text{GAN}} : X \rightarrow \mathbb{R}$	GAN discriminators <i>synthesized / real hand-only images</i>
$f^{\text{DAN}} : X \rightarrow X \times Y$	Domain adaptation network: $f^{\text{DAN}} = g^{\text{MR}} \circ g^{\text{FPE}}$ <i>hand-only / Ho images</i>

