3D Hand Pose Estimation via Regularized Graph Representation Learning

Anonymous Authors AAAI paper 6661

Abstract

This paper addresses the problem of 3D hand pose estimation 1 from a monocular RGB image. While previous methods have 2 shown great success, the structure of hands has not been fully 3 exploited, which is critical in pose estimation. To this end, we 4 propose a regularized graph representation learning under a 5 6 conditional adversarial learning framework for 3D hand pose estimation, aiming to capture structural inter-dependencies of 7 hand joints. In particular, we estimate an initial hand pose 8 from a parametric hand model as a prior of hand structure, 9 10 which regularizes the inference of the structural deformation in the prior pose for accurate graph representation learning 11 via residual graph convolution. To optimize the hand struc-12 13 ture further, we propose two bone-constrained loss functions, which characterize the morphable structure of hand poses ex-14 plicitly. Also, we introduce an adversarial learning framework 15 conditioned on the input image with a multi-source discrimi-16 nator, which imposes the structural constraints onto the distri-17 bution of generated 3D hand poses for anthropomorphically 18 valid hand poses. Extensive experiments demonstrate that our 19 model sets the new state-of-the-art in 3D hand pose estima-20 tion from a monocular image on five standard benchmarks. 21

1 Introduction

22

3D human hand pose estimation is a long-standing problem 23 in computer vision, which is critical for various applications 24 such as virtual reality and augmented reality (Hürst and van 25 Wezel 2011; Piumsomboon et al. 2013). Previous works at-26 tempt to estimate hand pose from depth images (Ge et al. 27 2016; Wu et al. 2018; Zhou et al. 2018; Ge et al. 2018) or 28 in multi-view setups (Panteleris and Argyros 2017; Zhang 29 et al. 2016a). However, due to the diversity and complexity 30 of hand shape, gesture, occlusion, etc., it still remains a chal-31 lenging problem despite years of studies (Rehg and Kanade 32 1994; Ying and Huang 2002; Ying, John, and Huang 2005; 33 Hui et al. 2017). 34

As RGB cameras are more widely accessible than depth sensors, recent works focus mostly on 3D hand pose estimation from a monocular RGB image and have shown their efficiency (Ge et al. 2019; Boukhayma, Bem, and Torr 2019; Baek, Kim, and Kim 2019; Cai et al. 2018; Zimmermann and Brox 2017a; Doosti et al. 2020). While some early works (Cai et al. 2018; Boukhayma, Bem, and Torr 2019)

Copyright © 2021, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.



Figure 1: The proposed method estimates 3D hand pose from a monocular image based on regularized graph representation learning. A parametric hand model generates a *prior pose*, which regularizes the learning of deformations in graph topology under a conditional adversarial learning framework.

did not explicitly exploit the structure of hands, some recent 42 methods (Ge et al. 2019; Doosti et al. 2020) have shown the 43 crucial role of hand structure in pose estimation, but may re-44 sort to an additional synthetic dataset. Also, unlike bodies 45 and faces that have obvious local characteristics (e.g., eyes 46 on a face), hands exhibit almost uniform appearance. Con-47 sequently, estimated hand poses from existing methods are 48 sometimes distorted and unnatural. 49

To fully exploit the structure of hands, we propose to rep-50 resent the irregular topology of 3D hand poses naturally on 51 graphs, and learn the graph representation regularized by a 52 prior pose from the monocular image input under a con-53 ditional generative adversarial learning framework, aiming 54 to capture the structural dependencies among hand joints. 55 Based on the Maximum a Posteriori estimation formulation 56 of inferring 3D hand pose, we first construct an initial hand 57 pose from a parametric hand model as a prior of hand struc-58 ture (prior pose), which captures the holistic topology of 59 hand structures, *i.e.*, the adjacency relations between joints. 60 Based on this prior pose, we represent the topology of hand 61 joints on a graph, where each joint is treated as a node and 62 each pair of adjacent nodes are connected. We further learn 63 the *deformation* in the prior pose to refine the hand struc-64 ture representation, by propagating information across ad-65 jacent nodes via residual graph convolution and conditional 66

on the input image. Moreover, while most existing works 67 (Boukhayma, Bem, and Torr 2019; Ge et al. 2019; Cai et al. 68 2018) deploy 3D Euclidean distance between joints as the 69 loss function for 3D annotation, we propose two bone loss 70 functions that constrain the length and orientation of each 71 bone connected by adjacent joints so as to preserve hand 72 structure explicitly. On the other hand, to address the chal-73 lenge of uniform appearance, we propose to train the net-74 work under an adversarial learning framework conditioned 75 on the input image, aiming to estimate the real distribution 76 of 3D hand poses. Besides, unlike some recent works (Ge 77 et al. 2019; Cai et al. 2018; Kulon et al. 2019), we estimate 78 3D hand poses without resorting to ground truth meshes or 79 depth maps, which is more suitable for datasets in the wild. 80

Specifically, given an input monocular image, our frame-81 work consists of a hand pose generator and a conditional 82 discriminator. The generator is composed of a MANO hand 83 model module (Romero, Tzionas, and Black 2017) that pro-84 vides an initial pose estimation as prior pose and a deforma-85 tion learning module regularized by the prior pose. In par-86 ticular, taking the prior pose and image features as input, 87 the deformation learning module learns the deformation in 88 the prior pose to further refine the hand structure, by our 89 designed residual graph convolution that leverages on the 90 recently proposed ResGCN (Li et al. 2019). Further, we de-91 sign a conditional multi-source discriminator that employs 92 hand poses, hand bones computed from poses as well as 93 the input image to distinguish the predicted 3D hand pose 94 from the ground-truth, leading to anthropomorphically valid 95 hand pose. Experimental results demonstrate that our model 96 achieves significant improvements over state-of-the-art ap-97 proaches on five standard benchmarks. 98

⁹⁹ To summarize, our main contributions include

• We propose regularized graph representation learning for 3D hand pose estimation from a monocular image, which fully exploits structural dependencies among hand joints.

• We learn the graph representation of hand poses by inferring structural deformation, which is regularized by an initial hand pose estimation from a parametric hand model.

• We introduce two bone-constrained loss functions, which optimize the estimation of hand structures by explicitly enforcing constrains on the topology of bones.

• We present a conditional adversarial learning framework to impose structural constraints onto the distribution of generated 3D hand poses, which is able to address the challenge of uniform appearance in hands.

2 Related Work

According to the input modalities, previous works on 3D
hand pose estimation can be classified into three categories:
1) 3D hand pose estimation from depth images; 2) 3D hand
pose estimation from multiple RGB images; 3) 3D hand
pose estimation from a monocular RGB image.

120 **2.1 Estimation from Depth Images**

114

121 Depth images contain rich 3D information for hand pose 122 estimation (Tang, Yu, and Kim 2013), which has shown promising accuracy (Yuan et al. 2018). There is a rich lit-123 erature on 3D hand pose estimation with depth images as 124 input (Ge et al. 2018, 2016; Fitzgibbon 2015; Choi 2016; 125 De, Fleet, and Paragios 2011; Khamis et al. 2015; Xiao 126 et al. 2015; Malik et al. 2018; Oberweger and Lepetit 2018). 127 Among them, some earlier works such as (De, Fleet, and 128 Paragios 2011; Khamis et al. 2015) are based on a de-129 formable hand model with an iterative optimization training 130 approach. Due to the effectiveness of deep learning, some 131 recent works like (Malik et al. 2018) leverage CNN to learn 132 the shape and pose parameters for a proposed model (LBS 133 hand model). 134

2.2 Estimation from Multiple Images

Multiple RGB images taken from different views also con-136 tain rich 3D information. Therefore, some works take multi-137 ple images as input to alleviate the occlusion problem (Cam-138 pos and Murray 2006; Oikonomidis, Kyriazis, and Argy-139 ros 2010; Sridhar et al. 2014). Campos et al. (Campos and 140 Murray 2006) propose a regression-based approach for hand 141 pose estimation, where they utilize multi-view images to 142 overcome the occlusion issue. Sridhar et al. (Sridhar et al. 143 2014) contribute a fundamentally extended generative track-144 ing algorithm based on an augmented implicit shape repre-145 sentation with multiple images as input. 146

135

147

2.3 Estimation from a Monocular Image

Compared with the aforementioned two categories, a 148 monocular RGB image is more accessible. Early works 149 (Athitsos and Sclaroff 2003; Rehg and Kanade 2002; 150 Stenger et al. 2006) propose complex model-fitting ap-151 proaches, which are based on dynamics and multiple hy-152 potheses and depend on restricted requirements. These 153 model-fitting approaches have proposed many hand mod-154 els, based on assembled geometric primitives (Oikonomidis, 155 Kyriazis, and Argyros 2011) or sphere meshes (Tkach, 156 Pauly, and Tagliasacchi 2016), etc. Our work deploys the 157 MANO hand model (Romero, Tzionas, and Black 2017) as 158 our prior, which models both hand shape and pose as well 159 as generates meshes. Nevertheless, these sophisticated ap-160 proaches suffer from low estimation accuracy. 161

With the advance of deep learning, many recent works 162 estimate 3D hand pose from a monocular RGB image us-163 ing neural networks (Ge et al. 2019; Boukhayma, Bem, and 164 Torr 2019; Baek, Kim, and Kim 2019; Cai et al. 2018; Zim-165 mermann and Brox 2017a; Yang and Yao 2019; Kulon et al. 166 2019). Among them, some recent works (Kulon et al. 2019; 167 Ge et al. 2019) directly reconstruct the 3D hand mesh and 168 then generate the 3D hand pose through a pose regressor. 169 Kulon et al. (Kulon et al. 2019) reconstruct the hand pose 170 based on an auto-encoder, which employs an encoder to ex-171 tract the latent code and feeds the latent code into the de-172 coder to reconstruct hand mesh. Ge et al. (Ge et al. 2019) 173 propose to estimate vertices of 3D meshes from GCNs (Kipf 174 and Welling 2017) in order to learn nonlinear variations in 175 hand shape. The latent feature of the input RGB image is 176 extracted via several networks and then fed into a GCN to 177 directly infer the 3D coordinates of mesh vertices. How-178 ever, since the accuracy of the output hand mesh is criti-179



Figure 2: Architecture of the proposed regularized graph representation learning under a conditional adversarial learning framework for 3D hand pose estimation.

cal for both methods, they need an extra dataset which provides ground truth hand meshes as supervision. Also, the
upsampling layer used in (Ge et al. 2019) to reconstruct the
hand mesh will cause a non-uniform distribution of vertices
in mesh, which influences the accuracy of hand pose.

In contrast, we take a prior pose estimated from a para-185 metric hand model as regularization for graph representation 186 learning over hand poses rather than directly reconstructing 187 hand poses from latent features. Besides, our method does 188 not require any additional supervision such as mesh super-189 vision (Ge et al. 2019; Kulon et al. 2019) or depth image 190 supervision (Ge et al. 2019; Cai et al. 2018). Hence, our 191 method is more suitable for datasets in the wild. Further, we 192 introduce conditional adversarial training for 3D hand pose 193 estimation, which enables learning a real distribution of 3D 194 hand poses. 195

3 Methodology

197 **3.1 Overview of the Proposed Approach**

196

We aim to infer 3D hand pose via regularized graph representation learning under an adversarial learning framework. The entire framework consists of a hand pose generator G and a conditional discriminator D, as illustrated in Fig. 2.

Given a monocular RGB image I as the input, the generator G includes two modules:

• The hand model module generates an initial estimation of 3D hand pose $\tilde{\mathbf{P}} \in \mathbb{R}^{N \times 3}$ with N joints (N = 21 in our experimental setting), which serves as a *prior pose* for the subsequent refinement. This module consists of a feature extractor and a parametric hand model.

• The deformation learning module infers the structural deformation in the prior pose $\tilde{\mathbf{P}}$ for regularized graph representation learning. Taking $\tilde{\mathbf{P}}$ and \mathbf{I} as the input, this module exploits the structural relationship among hand joints via residual graph convolution and outputs the deformation, leading to the refined pose $\hat{\mathbf{P}} \in \mathbb{R}^{N \times 3}$.

The multi-source discriminator \mathbb{D} imposes structural constraints onto the distribution of generated 3D hand poses conditioned on the input image, which distinguishes the 217 ground-truth 3D poses from the predicted ones. 218

3.2 The Proposed Hand Pose Generator \mathbb{G} 219

Given the observed input image I and ground truth hand pose P_{gt} , we formulate the training of hand pose estimation from a monocular image as a Maximum a Posteriori (MAP) estimation problem: 223

$$\hat{\mathbf{P}}_{MAP}(\mathbf{I}, \mathbf{P}_{gt}) = \operatorname*{argmax}_{\mathbf{P}} f(\mathbf{I}, \mathbf{P}_{gt} | \mathbf{P}) g(\mathbf{P}), \qquad (1)$$

where **P** denotes the hand pose to estimate. In (1), $g(\mathbf{P})$ represents the prior probability distribution of the hand pose, which provides the prior knowledge of **P**. $f(\mathbf{I}, \mathbf{P}_{gt} | \mathbf{P})$ denotes the likelihood function, which is the probability of obtaining the observed image **I** and ground truth hand pose \mathbf{P}_{gt} given the estimated hand pose **P**. 228

We define the likelihood function as an exponential function of the distance between the estimated pose and the ground truth pose/input image: 232

$$f(\mathbf{I}, \mathbf{P}_{gt}|\mathbf{P}) = \exp\{-d_1(\mathbf{P}_{gt}, \mathbf{P}) - d_2(\mathbf{I}, \mathbf{P})\}, \quad (2)$$

where $d_1(\cdot)$ is the distance metric between the estimated hand pose and the ground truth, and $d_2(\cdot)$ is the distance metric between the estimated hand pose and the input image. Regarding $g(\mathbf{P})$, it is a constant C after we acquire a prior pose from a parametric hand model. Hence, when we substitute (2) and $g(\mathbf{P}) = C$ into (1), take the logarithm and multiply by -1, we have 233

$$\min_{\mathbf{P}} d_1(\mathbf{P}_{\mathsf{gt}}, \mathbf{P}) + d_2(\mathbf{I}, \mathbf{P}). \tag{3}$$

 $d_1(\cdot)$ and $d_2(\cdot)$ will be discussed in Section 3.4 in detail.

Specifically, we employ a parametric hand model to provide the prior of **P**, and designate a Deformation Learning Module to learn the pose under the supervision of the ground-truth pose and input image. We discuss the two modules of the generator in detail as follows. 243

240



Figure 3: Illustration of the residual between the ground truth hand pose (marked in green) and the predicted one (marked in red). Each hand pose has 21 key joints. We denote a bone vector connecting two key joints i and j by $\mathbf{b}_{i,j}$, such as $\mathbf{b}_{5,6}$ in the figure.

The Hand Model Module Given an input monocular image, this module aims to generate an initial estimation of 3D hand pose $\tilde{\mathbf{P}}$ as a prior. A hand model is able to represent both hand shape and pose with a few parameters, which is thus a suitable prior for hand pose estimation.

We first predict parameters of the hand model. Specifi-251 cally, we crop and resize the input image to a salient region 252 of the hand, which is fed into the ResNet-50 network (He 253 et al. 2016) to extract features for the construction of the la-254 tent code z, *i.e.*, parameters of the hand model. Then, we 255 employ a modified MANO hand model (Romero, Tzionas, 256 and Black 2017), which is based on the SMPL model (Loper 257 et al. 2015) for human bodies. The MANO hand model is 258 a deformable hand mesh model with two vectors θ and β 259 contained in the latent code z as the input, which control 260 the pose and shape of the generated hand respectively. We 261 262 modify the default setting of $\{\theta, \beta\}$ from $\{10, 45\}$ to $\{10, 8\}$ for reduced computation complexity. Also, note that, while 263 Boukhayma et al. (Boukhayma, Bem, and Torr 2019) cre-264 ate a synthetic dataset to pre-train the ResNet-50 so as to 265 estimate parameters of MANO, we do not resort to any ex-266 tra dataset. The output of the MANO hand model includes a 267 hand pose $\mathcal{P}(\theta, \beta)$. 268

Additionally, we need to position the pose $\mathcal{P}(\theta, \beta)$ in a 269 camera coordinate system so as to acquire the 3D coordi-270 nates of each point in the hand pose. We project $\mathcal{P}(\theta, \beta)$ to 271 the 3D space via three parameters that model the camera co-272 ordinate system: 1) a 3D rotation parameter $\mathbf{c_r} \in \mathbb{R}^3$; 2) a 273 3D translation parameter $\mathbf{c_t} \in \mathbb{R}^3$; and 3) a scale param-274 eter $\mathbf{c}_{s} \in \mathbb{R}$. The camera parameters are estimated by the 275 aforementioned ResNet-50 network. 276

277 We formulate the complete hand model as:

$$\mathbf{P}(\theta, \beta, \mathbf{c_r}, \mathbf{c_t}, \mathbf{c_s}) = \mathbf{c_s} * R(\mathcal{P}(\theta, \beta), \mathbf{c_r}) + \mathbf{c_t}, \quad (4)$$

where R is a rotation function. The acquired initial estimation $\tilde{\mathbf{P}}$ serves as a prior pose for the subsequent deformation learning model.

The Deformation Learning Module This module aims at accurate graph representation learning for hand pose estima-



Figure 4: Architecture of the deformation learning module in our generator.

tion, which is conditional on the prior and under the supervision of the input image and ground truth pose as in (1). In particular, conditioned on the prior $\tilde{\mathbf{P}}$, we learn the structural *deformation* in $\tilde{\mathbf{P}}$ instead of the holistic hand pose. 286

We first construct an unweighted graph over \mathbf{P} , where the 287 irregularly sampled key points (*i.e.*, joints) on the hand are 288 projected onto nodes. The graph signal on each node is the 289 concatenation of the global feature vector of the input im-290 age and the 3-dimensional coordinate vector of each joint in 291 the input prior pose. Nodes are connected if they represent 292 adjacent key points of the hand, where the adjacency rela-293 tions follow the human hand structure as presented in Fig. 3, 294 leading to an adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$. 295

Based on the graph representation **A**, the finally refined pose is

$$\hat{\mathbf{P}} = \tilde{\mathbf{P}} + \operatorname{GCN}(\tilde{\mathbf{P}} \oplus \mathbf{F}, \mathbf{A}), \tag{5}$$

296

297

where $\mathbf{F} \in \mathbb{R}^{N \times F}$ denotes the *F*-dimensional global feature vector of the image repeated *N* times, and \oplus denotes the feature-wise concatenation operation. $\operatorname{GCN}(\tilde{\mathbf{P}} \oplus \mathbf{F}, \mathbf{A})$ the ground truth. The sum of the prior pose $\tilde{\mathbf{P}}$ and its deformation thus leads to the refined hand pose. 303

Specifically, we first extract features from the RGB image 304 I to facilitate the pose refinement. We follow the ResNet-50 305 architecture (He et al. 2016) to extract 2D features F from 306 I. Next, we learn the structural deformation in the hand pose 307 via residual graph convolution. Leveraging on the idea of 308 the recent ResGCN (Li et al. 2019) that shows graph resid-309 ual learning enables deeper graph convolution networks and 310 better feature learning, we design a Graph Res-block to learn 311 the deformation of the prior pose. Specifically, we employ 312 the efficient GCN (Kipf and Welling 2017) as the basic unit 313 of the Graph Res-block, which essentially propagates infor-314 mation across adjacent nodes to learn higher-level features. 315 Each Graph Res-block consists of two GCN layers as well 316 as two normalization layers that enable higher learning rate 317 without vanishing or exploding gradients. Furthermore, we 318 introduce residual skip connections for all the Graph Res-319 blocks in order to accelerate the speed of convergence and 320 avoid the gradient vanishment. 321

Let \mathbf{X}^{l} denote the input of the *l*-th Graph Res-block, then the output of the *l*-th Graph Res-block takes the form

$$\mathbf{X}^{l+1} = N\left(g(N(g(\mathbf{X}^l, \mathbf{A})), \mathbf{A})\right) + \operatorname{skip}(\mathbf{X}^l), \quad (6)$$



Figure 5: Architecture of the conditional discriminator.

where $g(\cdot)$ represents a single GCN layer as in (Kipf and Welling 2017), $N(\cdot)$ represents a single normalization layer, and skip(\cdot) denotes the skip connection which is a GCN layer to match the dimension of the two terms in (6). We then stack several layers of Graph Res-blocks to learn the deformation of the prior pose, as demonstrated in Fig. 4.

330 3.3 The Proposed Conditional Discriminator $\mathbb D$

In the adversarial training stage, while we learn the generator to predict hand poses which are indistinguishable to the discriminator, the discriminator attempts to distinguish real samples from fake ones, *i.e.*, the predicted hand poses. In particular, given the input image **I**, we designate a *conditional* discriminator conditioned on **I**.

A simple architecture of a discriminator is a fully-337 connected (FC) network with the hand pose as input, which 338 however has two shortcomings: 1) the relation between the 339 RGB image and inferred hand pose is neglected; 2) struc-340 tural properties of the hand pose are not taken into account 341 explicitly. Instead, inspired by the multi-source architecture 342 in (Yang et al. 2018), we design a conditional multi-source 343 discriminator with three inputs to address the aforemen-344 tioned issues. As illustrated in Fig. 5, the inputs include: 1) 345 features of the input monocular image; 2) features of the re-346 fined hand pose $\hat{\mathbf{P}}$ or the ground truth pose \mathbf{P}_{gt} ; 3) features 347 of bones via the KCS layer as in (Wandt, Ackermann, and 348 Rosenhahn 2017), which computes the bone matrix from $\hat{\mathbf{P}}$ 349 or \mathbf{P}_{gt} via a simple matrix multiplication. The bone features 350 contain prominent structural information such as the length 351 and direction of bones, thus characterizing the hand struc-352 ture accurately. 353

The loss function of the conditional discriminator follows the definition of the Wasserstein loss (Arjovsky, Chintala, and Bottou 2017) conditioned on the input image I:

$$\mathcal{L}_{\text{Wass}} = -\mathbb{E}_{\mathbf{P}_{\text{gt}} \sim p_{data}(\mathbf{P}_{\text{gt}})} \mathbb{D}(\mathbf{P}_{\text{gt}} | \mathbf{I}) + \mathbb{E}_{\hat{\mathbf{P}} \sim p(\hat{\mathbf{P}})} \mathbb{D}(\hat{\mathbf{P}} | \mathbf{I}),$$
(7)

where \mathbb{D} takes the generated (fake) pose $\hat{\mathbf{P}}$ and ground-truth pose \mathbf{P}_{gt} as input, \mathbf{P}_{gt} is a sample following the ground-truth pose distribution $p_{data}(\mathbf{P}_{gt})$ and $\hat{\mathbf{P}}$ is a sample from the refined pose distribution $p(\hat{\mathbf{P}})$.

Specifically, we employ a CNN to extract features of the input monocular image, a GCN to learn the representation of the refined pose or the ground truth pose, and one FC layer to capture the features of bone structures computed from the hand pose. Besides, the architecture of our multi-source discriminator is based on SNGAN (Miyato et al. 2018) with spectral normalization layers. 367

3.4 The Proposed Bone-Constrained Loss 368 Functions 369

As presented in (3), we have two types of loss functions for 370 the MAP estimation of hand pose. We employ the commonly 371 adopted average Euclidean distance in the coordinates of 372 joints of 3D hand pose \mathcal{L}_{pose} (Ge et al. 2019) as well as 373 two proposed bone-constrained metrics as $d_1(\cdot)$ to measure 374 the distortion of the estimated 3D hand pose compared to 375 the ground truth, and apply the commonly used average Eu-376 clidean distance in the coordinates of joints of projected 2D 377 hand pose $\mathcal{L}_{\text{proj}}$ (Ge et al. 2019) as $d_2(\cdot)$ to measure the dis-378 tance between the estimation and the 2D image. 379

Since \mathcal{L}_{pose} and \mathcal{L}_{proj} cannot capture the structural properties of hand pose explicitly, we propose two novel boneconstrained loss functions to characterize the length and direction of each bone.

As illustrated in Fig. 3, we first define a bone vector $\mathbf{b}_{i,j} \in \mathbb{R}^{3 \times 1}$ between hand joint i and j as

$$\mathbf{b}_{i,j} = \mathbf{j}_i - \mathbf{j}_j,\tag{8}$$

where $\mathbf{j}_i, \mathbf{j}_j \in \mathbb{R}^{3 \times 1}$ are the coordinates of joint i and j respectively.

The first bone-constrained loss \mathcal{L}_{len} quantifies the distance in *bone length* between the ground truth hand and its estimate, which we define as 390

$$\mathcal{L}_{\text{len}} = \sum_{i,j} \left| ||\mathbf{b}_{i,j}||_2 - ||\hat{\mathbf{b}}_{i,j}||_2 \right|,$$
(9)

where $\mathbf{b}_{i,j}$ and $\hat{\mathbf{b}}_{i,j}$ are the bone vectors of the ground truth and the predicted bone respectively. 392

The second bone-constrained loss \mathcal{L}_{dir} measures the deviation in the *direction of bones*: 393

$$\mathcal{L}_{\text{dir}} = \sum_{i,j} \left| \left| \mathbf{b}_{i,j} / || \mathbf{b}_{i,j} ||_2 - \hat{\mathbf{b}}_{i,j} / || \hat{\mathbf{b}}_{i,j} ||_2 \right| \right|_2.$$
(10)

This is motivated by the fact that small loss in joints some-395 times may not reflect large distortion in hand pose. Taking 396 joints \mathbf{j}_5 and \mathbf{j}_6 in Fig. 3 as an example, the distance between 397 the ground truth joints and predicted ones is trivial. How-398 ever, it is obvious that the orientation of the predicted bone 399 $\hat{\mathbf{b}}_{5,6}$ significantly deviates from the ground truth $\mathbf{b}_{5,6}$. This 400 distortion in hand structure is well captured by our proposed 401 loss in the bone direction \mathcal{L}_{dir} . 402

Besides, as we adopt the framework of adversarial learning, we also introduce the Wasserstein loss \mathcal{L}_{Wass} in (7) into the loss function for adversarial training. Hence, the overall loss function \mathcal{L} is 406

$$\mathcal{L} = \mathcal{L}_{\text{pose}} + \lambda_{\text{proj}} \mathcal{L}_{\text{proj}} + \lambda_{\text{len}} \mathcal{L}_{\text{len}} + \lambda_{\text{dir}} \mathcal{L}_{\text{dir}} + \lambda_{\text{Wass}} \mathcal{L}_{\text{Wass}},$$
(11)

where λ_{proj} , λ_{len} , λ_{dir} and λ_{Wass} are hyperparameters for the 407 trade-off among these losses. In accordance with (3), $d_1 = 408 \mathcal{L}_{\text{pose}} + \lambda_{\text{len}} \mathcal{L}_{\text{len}} + \lambda_{\text{dir}} \mathcal{L}_{\text{dir}}$, and $d_2 = \lambda_{\text{proj}} \mathcal{L}_{\text{proj}}$.

387 388 389

395

Stage	hand model	deformation	discriminator S	STB	RHD	EGODEXTER
Ι	 ✓ 		2	24.15	83.37	52.32
Π	\checkmark	\checkmark		5.12	15.84	43.26
III	\checkmark	\checkmark	√ .	3.97	12.40	34.98

Table 1: The performance of different stages in our model on three datasets (measured in 3D Euclidean distance (mm)).

	GCN Deformation	FC Deformation	Discriminator	STB	RHD	EGODEXTER
1		✓		15.11	37.59	52.34
2	\checkmark			5.12	15.84	40.12
3		\checkmark	\checkmark	10.23	25.15	44.23
4	\checkmark		\checkmark	3.97	12.40	34.98

Table 2: Ablation studies on the Deformation Learning Module, with comparison between the Deformation Learning Module and the simple FC Refinement Module in 3D Euclidean distance (mm).

4 Experimental Results

411 **4.1 Datasets and Metrics**

410

Datasets We evaluate our approach on five public datasets: 412 Stereo Hand Pose Tracking Benchmark (STB) (Zhang 413 et al. 2016b), the Rendered Hand Pose Dataset (RHD) 414 415 (Zimmermann and Brox 2017b), EGODEXTER (Mueller 416 et al. 2017), MPII+NZSL (Simon et al. 2017) and DEX-TER+OBJECT (Mueller et al. 2017). We use STB, RHD 417 and EGODEXTER for ablation studies and compare with 418 state-of-the-art methods on all the five datasets. 419

STB is a real-world dataset with image resolution of 420 640×320 . Following (Zimmermann and Brox 2017b), we 421 422 split the 18,000 images into 15,000 training samples and 3,000 test samples. Besides, to make the definition of joints 423 consistent, we move the location of the root joint from the 424 palm center to the wrist following (Ge et al. 2019). RHD is 425 a more challenging synthetic dataset with image resolution 426 of 320×320 , which is built upon 20 different characters 427 428 performing 39 actions. MPII+NZSL is a real-world dataset 429 containing images from YouTube videos. This dataset only provides 2D annotations. DEXTER+OBJECT dataset shows 430 interactions of an actor's hand with a cuboid object from 431 a third person view. EGODEXTER dataset displays a hand 432 from an egocentric view interacting with various objects. 433

Metrics We evaluate the performance of 3D hand pose
estimation with two metrics: (i) pose error, which takes the
average location error in Euclidean distance between the estimated 3D joints and the ground truth; (ii) percentage of
correct key points (PCK), which is the percentage of correct key points whose error in Euclidean distance is below a
threshold.

441 **4.2 Implementation Details**

In our experiments, we first pretrain the hand model module
and then train the entire network end-to-end. In particular,
the training process can be divided into three stages.

Stage I. We pretrain the hand model module, which is
randomly initialized and trained for 100 epochs using the
Adam optimizer with learning rate 0.001. Then, we freeze
the parameters of this stage to evaluate the effectiveness of
the deformation learning module.

450 **Stage II.** We train the generator \mathbb{G} end-to-end without the 451 discriminator \mathbb{D} . In \mathbb{G} , the hand model module is initialized



Figure 6: Qualitative results of different stages in our model.

with the trained model in the first stage and the deformation 452 learning module is randomly initialized. G is then trained 453 with 100 epochs using the Adam optimizer with learning 454 rate 0.0001. 455

Stage III. We adopt the framework of SNGAN (Miyato456et al. 2018) for the conditional adversarial training, and train457our model end-to-end. G and D are trained with 100 epochs458using the Adam optimizer with learning rate 0.0001.459

Regarding the hyper-parameters in (11), we set $\lambda_{\text{len}} = 460$ $0.01, \lambda_{\text{dir}} = 0.1, \lambda_{\text{proj}} = 0.1, \lambda_{\text{Wass}} = 0.01.$

462

4.3 Ablation Studies

We perform ablation studies on the performance of different stages, the deformation learning module, the discriminator and loss functions. Due to the page limit, we present all the results in 3D Euclidean distance (mm). Please refer to the supplementary material for the results measured in 3D PCK.

On different stages. We present the results of three train-468 ing stages in average 3D Euclidean distance, as listed in 469 Tab. 1. The performance of Stage II significantly outper-470 forms Stage I, which demonstrates that the proposed defor-471 mation learning module plays the most critical role in our 472 model. The adversarial training scheme (Stage III) further 473 improves the result, by learning a real distribution of the 474 3D hand pose. We also show visual results of our method 475 at different stages in Fig. 6. We see that Stage I estimates 476 a coarse hand pose from the MANO hand model as a prior 477 pose, while Stage II refines the structure of the prior pose 478 significantly. Finally, Stage III generates more realistic hand 479 poses via conditional adversarial learning. 480

On the deformation learning module. We compare the 481 deformation learning module with a simple fully-connected 482 deformation learning module (FC Deformation Module) to 483 refine the prior pose. We train the deformation learning mod-484 ules in different experimental settings: 1) without our dis-485 criminator, *i.e.*, without adversarial learning; and 2) with our 486 discriminator. As presented in Tab. 2, the GCN deformation 487 learning module leads to significant gain over the simple FC 488 deformation module on both datasets in different settings, 489 thus validating the superiority of the proposed deformation 490 learning module. 491

On the conditional discriminator. We compare with a single-source discriminator which only takes the 3D hand pose as the input. As presented in Tab. 3, the multi-source discriminator outperforms the single-source one on both datasets, which gives credits to exploring the structure of hand bones and the relation between the image and pose. 497

De	formation Learning	Multi-source	Single-source	STB	RHD	EGODEXTER
$\begin{bmatrix} 1\\2 \end{bmatrix}$	√ √	\checkmark	✓	3.97 4.54	12.40 15.10	34.98 37.46

Table 3: Ablation studies on the discriminator (3D Euclidean distance (mm)).

	C + C	C.	\mathcal{L}_{dir}	STB			RHD		
L _{pose} +	$\mathcal{L}_{pose} + \mathcal{L}_{proj}$	Llen		Stage I	Stage II	Stage III	Stage I	Stage II	Stage III
1	✓			32.75	9.11	5.35	99.24	25.96	15.07
2	\checkmark	\checkmark		30.32	8.00	5.02	95.19	22.96	14.76
3	~		\checkmark	27.65	6.91	5.00	89.76	21.63	14.01
4	√	\checkmark	\checkmark	24.15	5.12	3.97	83.37	15.84	12.40

Table 4: Ablation studies on the proposed bone-constrained loss functions at three stages.

On loss functions. We also evaluate the proposed bone-498 constrained loss functions \mathcal{L}_{len} and \mathcal{L}_{dir} separately. We train 499 the network with different combinations of loss functions on 500 the STB and RHD datasets in three stages respectively. As 501 reported in Tab. 4, the network trained with our proposed 502 bone-constrained loss functions performs better in all the 503 three stages on both datasets. We also notice that \mathcal{L}_{dir} plays 504 a more significant role compared to \mathcal{L}_{len} . This gives cred-505 its to the constraint on the orientation of bones that explic-506 itly takes structural properties of hands into consideration. 507 Further, we demonstrate the visual comparison of estimated 508 poses with and without bone-constrained losses in Fig. 7. 509 The estimated pose may have unnatural distortion in the di-510 rection of bones in the absence of the bone-constrained loss 511 functions, e.g., the little finger in the first row and the thumb 512 in the second row. In contrast, our results exhibit natural and 513 accurate structure in the orientation of bones with the pro-514 posed bone constraints enforced. 515

516 4.4 Experimental Results

	STB	RHD	MPII+ZNSL(px)	Dexter+Object	EgoDexter
(Ge et al. 2019)	6.37	15.33	-	-	-
(Boukhayma, Bem, and Torr 2019)	9.76	-	18.95	25.53	45.33
(Spurr et al. 2018)	8.56	19.73	-	40.20	56.92
(Zimmermann and Brox 2017b)	-	-	59.40	34.75	52.77
Ours	3.97	12.40	9.87	16.12	34.98

Table 5: Comparison with state-of-the-art methods on the five datasets. Note that MPII+ZNSL only provides 2D annotation, thus we employ the 2D distance (px) metric on this dataset.

We compare our method with competitive 3D hand pose 517 estimation approaches on the five datasets. On the relatively 518 simple STB dataset, as shown in Fig. 8, we compare with lat-519 est methods (Cai et al. 2018; Iqbal et al. 2018; Boukhayma, 520 Bem, and Torr 2019; Ge et al. 2019; Mueller et al. 2018a; 521 Zimmermann and Brox 2017a; Mueller et al. 2018b). Our 522 paradigm outperforms the state-of-the-art (Ge et al. 2019), 523 which closely reaches the upper bound 1.0 of 3D PCK at all 524 the error thresholds. Also, we list the results in 3D Euclidean 525 distance for comparison with the state-of-the-arts in Tab. 5. 526 Compared to these works which directly reconstruct the 3D 527 hand pose (Ge et al. 2019; Boukhayma, Bem, and Torr 2019; 528 Cai et al. 2018), our method performs much better mainly 529 due to the proposed regularized graph representation learn-530 ing and conditional adversarial learning. 531

532 Moreover, we demonstrate some qualitative results of our



Figure 7: Qualitative results for the evaluation of the proposed bone-constrained loss functions.



Figure 8: Comparison with state-of-the-art methods on the STB dataset.



Figure 9: Qualitative results of our proposed network on the STB dataset. The 2D pose is projected from the 3D pose.

3D hand pose estimation in Fig. 9. The generated poses are accurate and natural even in case of severe self-occlusions, as shown in the first three columns of Fig. 9. This validates the effectiveness of the proposed paradigm. We show the qualitative results and PCK results on the other four datasets in the supplementary material.

5 Conclusion

539

In this paper, we propose regularized graph representation 540 learning under a conditional adversarial learning framework 541 for 3D hand pose estimation from a monocular image. Based 542 on the MAP estimation formulation, we take an initial es-543 timation of hand pose as prior pose, and further learn the 544 structural deformation in the prior pose via residual graph 545 convolution. Also, we propose two bone-constrained loss 546 functions to enforce constraints on the bone structures ex-547 plicitly. Extensive experiments demonstrate the superiority 548 of the proposed method. 549

References

551 Arjovsky, M.; Chintala, S.; and Bottou, L. 2017. Wasser-

stein Generative Adversarial Networks. In Precup, D.; and

Teh, Y. W., eds., Proceedings of the 34th International Con-

ference on Machine Learning, volume 70 of *Proceedings of Machine Learning Research*, 214–223. International Con-

vention Centre, Sydney, Australia: PMLR.

557 Athitsos, V.; and Sclaroff, S. 2003. Estimating 3D Hand

Pose From a Cluttered Image. In *IEEE Computer Society Conference on Computer Vision & Pattern Recognition.*

560 Baek, S.; Kim, K. I.; and Kim, T.-K. 2019. Pushing the En-

velope for RGB-Based Dense 3D Hand Pose Estimation via

562 Neural Rendering. In *The IEEE Conference on Computer*

563 Vision and Pattern Recognition (CVPR).

⁵⁶⁴ Boukhayma, A.; Bem, R. d.; and Torr, P. H. 2019. 3D Hand

Shape and Pose From Images in the Wild. In *The IEEE Conference on Computer Vision and Pattern Recognition* (CVPR).

Cai, Y.; Ge, L.; Cai, J.; and Yuan, J. 2018. Weaklysupervised 3D Hand Pose Estimation from Monocular RGB
Images. In *The European Conference on Computer Vision*

571 (ECCV).

⁵⁷² Campos, T. E. D.; and Murray, D. W. 2006. Regression-⁵⁷³ based Hand Pose Estimation from Multiple Cameras. In

⁵⁷³ based Hand Pose Estimation from Multiple Cameras. In
 ⁵⁷⁴ *IEEE Computer Society Conference on Computer Vision & Pattern Recognition.*

576 Choi, C. 2016. DeepHand: Robust Hand Pose Estimation

577 by Completing a Matrix Imputed with Deep Features. In 578 *Computer Vision & Pattern Recognition*.

De, L. G. M.; Fleet, D. J.; and Paragios, N. 2011. ModelBased 3D Hand Pose Estimation from Monocular Video. *IEEE Trans Pattern Anal Mach Intell* 33(9): 1793–1805.

582 Doosti, B.; Naha, S.; Mirbagheri, M.; and Crandall, D. J.

583 2020. HOPE-Net: A Graph-Based Model for Hand-Object

Pose Estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).*

Fitzgibbon, A. 2015. Accurate, Robust, and Flexible Real-time Hand Tracking. *Inproceedings* 3633–3642.

588 Ge, L.; Cai, Y.; Weng, J.; and Yuan, J. 2018. Hand PointNet:

⁵⁸⁹ 3D Hand Pose Estimation Using Point Sets. 8417–8426. ⁵⁹⁰ doi:10.1109/CVPR.2018.00878.

Ge, L.; Liang, H.; Yuan, J.; and Thalmann, D. 2016. Robust 3D Hand Pose Estimation in Single Depth Images:
From Single-View CNN to Multi-View CNNs. In *The IEEE*

594 Conference on Computer Vision and Pattern Recognition 595 (CVPR).

⁵⁹⁶ Ge, L.; Ren, Z.; Li, Y.; Xue, Z.; Wang, Y.; Cai, J.; and Yuan,

J. 2019. 3D Hand Shape and Pose Estimation From a Single
 RGB Image. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual

Learning for Image Recognition. In *The IEEE Conference* on *Computer Vision and Pattern Recognition (CVPR)*. Hui, L.; Yuan, J.; Lee, J.; Ge, L.; and Thalmann, D. 2017. Hough Forest With Optimized Leaves for Global Hand Pose Estimation With Arbitrary Postures. *IEEE Transactions on Cybernetics* PP(99): 1–15.

Hürst, W.; and van Wezel, C. 2011. Gesture-based interaction via finger tracking for mobile augmented reality. *Multimedia Tools and Applications* 62: 233–258. 609

Iqbal, U.; Molchanov, P.; Breuel Juergen Gall, T.; and Kautz,610J. 2018. Hand Pose Estimation via Latent 2.5D Heatmap Regression. In The European Conference on Computer Vision611(ECCV).613

Khamis, S.; Taylor, J.; Shotton, J.; Keskin, C.; Izadi, S.; and Fitzgibbon, A. 2015. Learning an efficient model of hand shape variation from depth images. In *IEEE Conference on Computer Vision & Pattern Recognition.* 617

Kipf, T. N.; and Welling, M. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations (ICLR).* 620

Kulon, D.; Wang, H.; Güler, R. A.; Bronstein, M. M.; and Zafeiriou, S. 2019. Single Image 3D Hand Reconstruction with Mesh Convolutions. In *BMVC*. 623

Li, G.; Muller, M.; Thabet, A.; and Ghanem, B. 2019. Deep-GCNs: Can GCNs Go As Deep As CNNs? In *The IEEE International Conference on Computer Vision (ICCV).* 626

Loper, M.; Mahmood, N.; Romero, J.; Pons-Moll, G.; and Black, M. J. 2015. SMPL: A Skinned Multi-person Linear Model. *ACM Trans. Graph.* 34(6): 248:1–248:16. ISSN 0730-0301. doi:10.1145/2816795.2818013. URL http://doi. acm.org/10.1145/2816795.2818013. 631

Malik, J.; Elhayek, A.; Nunnari, F.; Varanasi, K.; and Stricker, D. 2018. DeepHPS: End-to-end Estimation of 3D Hand Pose and Shape by Learning from Synthetic Depth. In 2018 International Conference on 3D Vision (3DV). 635

Miyato, T.; Kataoka, T.; Koyama, M.; and Yoshida, Y. 2018. 636 Spectral Normalization for Generative Adversarial Networks. In *International Conference on Learning Represen-* 638 *tations*. URL https://openreview.net/forum?id=B1QRgziT-. 639

Mueller, F.; Bernard, F.; Sotnychenko, O.; Mehta, D.; Sridhar, S.; Casas, D.; and Theobalt, C. 2018a. GANerated Hands for Real-Time 3D Hand Tracking from Monocular RGB. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*. URL https://handtracker.mpiinf.mpg.de/projects/GANeratedHands/. 645

Mueller, F.; Bernard, F.; Sotnychenko, O.; Mehta, D.; Srid-
har, S.; Casas, D.; and Theobalt, C. 2018b. GANerated
Hands for Real-Time 3D Hand Tracking From Monocular
RGB. In The IEEE Conference on Computer Vision and
Pattern Recognition (CVPR).648

Mueller, F.; Mehta, D.; Sotnychenko, O.; Sridhar, S.; Casas, 651 D.; and Theobalt, C. 2017. Real-time Hand Tracking under Occlusion from an Egocentric RGB-D Sensor. In *Proceedings of International Conference on Computer Vision* 654 *(ICCV)*. URL https://handtracker.mpi-inf.mpg.de/projects/ OccludedHands/. 656

550

- 657 Oberweger, M.; and Lepetit, V. 2018. DeepPrior++: Improv-
- 658 ing Fast and Accurate 3D Hand Pose Estimation. In 2017
- 659 IEEE International Conference on Computer Vision Work-
- 660 *shops (ICCVW)*.
- Oikonomidis, I.; Kyriazis, N.; and Argyros, A. 2011. Effi-
- cient model-based 3D tracking of hand articulations using
- 663 Kinect. volume 1. doi:10.5244/C.25.101.
- Oikonomidis, I.; Kyriazis, N.; and Argyros, A. A. 2010.
- Markerless and Efficient 26-DOF Hand Pose Recovery. In
 Computer Vision-accv -asian Conference on Computer Vi- sion.
- 668 Panteleris, P.; and Argyros, A. A. 2017. Back to RGB:
- 669 3D tracking of hands and hand-object interactions based
- on short-baseline stereo. *CoRR* abs/1705.05301. URL http://arxiv.org/abs/1705.05301.
- 672 Piumsomboon, T.; Clark, A.; Billinghurst, M.; and Cock-
- 673 burn, A. 2013. User-Defined Gestures for Augmented Re-
- 674 ality. In Kotzé, P.; Marsden, G.; Lindgaard, G.; Wesson,
- ⁶⁷⁵ J.; and Winckler, M., eds., *Human-Computer Interaction* ⁶⁷⁶ – *INTERACT 2013*, 282–299. Berlin, Heidelberg: Springer
- 677 Berlin Heidelberg.
- Rehg, J. M.; and Kanade, T. 1994. Visual tracking of high
 DOF articulated structures: An application to human hand
 tracking. In *Proc of Third European Conference on Com- puter Vision.*
- Rehg, J. M.; and Kanade, T. 2002. DigitEyes: vision-based
 hand tracking for human-computer interaction. In *IEEE Workshop on Motion of Non-rigid & Articulated Objects*.
- Romero, J.; Tzionas, D.; and Black, M. J. 2017. Em-
- bodied Hands: Modeling and Capturing Hands and Bodies Together. *ACM Trans. Graph.* 36(6): 245:1–245:17.
- 688 ISSN 0730-0301. doi:10.1145/3130800.3130883. URL
- 689 http://doi.acm.org/10.1145/3130800.3130883.
- 690 Simon, T.; Joo, H.; Matthews, I.; and Sheikh, Y. 2017. Hand
- 691 Keypoint Detection in Single Images Using Multiview Boot-
- strapping. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*
- Spurr, A.; Song, J.; Park, S.; and Hilliges, O. 2018. Crossmodal Deep Variational Hand Pose Estimation. *CoRR*abs/1803.11404. URL http://arxiv.org/abs/1803.11404.
- ⁶⁹⁷ Sridhar, S.; Rhodin, H.; Seidel, H. P.; Oulasvirta, A.; and
- Sridhar, S.; Rhodin, H.; Seidel, H. P.; Oulasvirta, A.; and
 Theobalt, C. 2014. Real-time hand tracking using a sum of
- anisotropic gaussians model. In *International Conference*on 3D Vision.
- Stenger, B.; Thayananthan, A.; Torr, P. H. S.; and Cipolla,
 R. 2006. Model-based hand tracking using a hierarchical Bayesian filter. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 28(9): 1372–1384. doi:10.1109/
 TPAMI.2006.189.
- 706 Tang, D.; Yu, T. H.; and Kim, T. K. 2013. Real-Time Artic-
- ⁷⁰⁷ ulated Hand Pose Estimation Using Semi-supervised Trans ⁷⁰⁸ ductive Regression Forests. In *IEEE International Confer-*
- ence on Computer Vision.
- Tkach, A.; Pauly, M.; and Tagliasacchi, A. 2016. Spheremeshes for Real-time Hand Modeling and Tracking. *ACM*

- *Trans. Graph.* 35(6): 222:1–222:11. ISSN 0730-0301. 712 doi:10.1145/2980179.2980226. URL http://doi.acm.org/10. 713 1145/2980179.2980226. 714
- Wandt, B.; Ackermann, H.; and Rosenhahn, B. 2017. A 715 Kinematic Chain Space for Monocular Motion Capture . 716
- Wu, X.; Finnegan, D.; O'Neill, E.; and Yang, Y.-L. 2018.717HandMap: Robust Hand Pose Estimation via Intermediate718Dense Guidance Map Supervision. In The European Con-
ference on Computer Vision (ECCV).720
- Xiao, S.; Wei, Y.; Shuang, L.; Tang, X.; and Jian, S. 2015. 721 Cascaded hand pose regression. In *IEEE Conference on 722 Computer Vision & Pattern Recognition.* 723
- Yang, L.; and Yao, A. 2019. Disentangling Latent Hands for Image Synthesis and Pose Estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).* 726
- Yang, W.; Ouyang, W.; Wang, X.; Ren, J.; Li, H.; and Wang, X. 2018. 3D Human Pose Estimation in the Wild by Adversarial Learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).* 730
- Ying, W.; and Huang, T. S. 2002. Hand modeling, analysis 731 and recognition. *IEEE Signal Processing Magazine* 18(3): 732 51–60. 733
- Ying, W.; John, L.; and Huang, T. S. 2005. Analyzing and capturing articulated hand motion in image sequences. *IEEE Transactions on Pattern Analysis & Machine Intelligence* 27(12): 1910–1922. 737
- Yuan, S.; Garcia-Hernando, G.; Stenger, B.; Moon, G.; 738 Chang, J.; Lee, K.; Molchanov, P.; Kautz, J.; Honari, S.; 739 Ge, L.; Yuan, J.; Chen, X.; Wang, G.; Yang, F.; Akiyama, 740 K.; Wu, Y.; Wan, Q.; Madadi, M.; Escalera, S.; and Kim, 741 T.-K. 2018. Depth-Based 3D Hand Pose Estimation: From 742 Current Achievements to Future Goals. 2636–2645. doi: 743 10.1109/CVPR.2018.00279. 744
- Zhang, J.; Jiao, J.; Chen, M.; Qu, L.; Xu, X.; and Yang, Q. 745 2016a. 3D Hand Pose Tracking and Estimation Using Stereo 746 Matching . 747
- Zhang, J.; Jiao, J.; Chen, M.; Qu, L.; Xu, X.; and Yang, Q. 748 2016b. 3D Hand Pose Tracking and Estimation Using Stereo 749 Matching. *ArXiv* abs/1610.07214. 750
- Zhou, Y.; Lu, J.; Du, K.; Lin, X.; Sun, Y.; and Ma, X. 2018. 751 HBE: Hand Branch Ensemble Network for Real-time 3D 752 Hand Pose Estimation. In *The European Conference on* 753 *Computer Vision (ECCV)*. 754
- Zimmermann, C.; and Brox, T. 2017a. Learning to Estimate 755 3D Hand Pose from Single RGB Images. 4913–4921. doi: 756 10.1109/ICCV.2017.525. 757
- Zimmermann, C.; and Brox, T. 2017b. Learning to Estimate 758 3D Hand Pose From Single RGB Images. In *The IEEE International Conference on Computer Vision (ICCV).* 760