# 3D Hand Pose Estimation via Regularized Graph Representation Learning [*],[**]

Yiming He[1] and Wei Hu[1]

Wangxuan Institute of Computer Technology, Peking University

**Abstract.** This paper addresses the problem of 3D hand pose estimation from a monocular RGB image. While previous methods have shown great success, the structure of hands has not been fully exploited, which is critical in pose estimation. To this end, we propose a regularized graph representation learning under a conditional adversarial learning framework for 3D hand pose estimation, aiming to capture structural interdependencies of hand joints. In particular, we estimate an initial hand pose from a parametric hand model as a prior of hand structure, which regularizes the inference of the structural deformation in the prior pose for accurate graph representation learning via residual graph convolution. To optimize the hand structure further, we propose two bone-constrained loss functions, which characterize the morphable structure of hand poses explicitly. Also, we introduce an adversarial learning framework conditioned on the input image with a multi-source discriminator, which imposes the structural constraints onto the distribution of generated 3D hand poses for anthropomorphically valid hand poses. Extensive experiments demonstrate that our model sets the new state-of-the-art in 3D hand pose estimation from a monocular image on five standard benchmarks.

**Keywords:** 3D hand pose estimation· graph refinement· prior pose· adversarial learning· bone-constrained loss.

## 1 Introduction

3D human hand pose estimation is a long-standing problem in computer vision, which is critical for various applications such as virtual reality and augmented reality [15,25]. Previous works attempt to estimate hand pose from depth images [11,10] or in multi-view setups [24,33]. However, due to the diversity and complexity of hand shape, gesture, occlusion, *etc.*, it still remains a challenging problem despite years of studies [14].

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 [12,4,3,5,8]. While some early works [5,4] did not explicitly exploit the structure of hands, some recent methods [12,8] have shown the

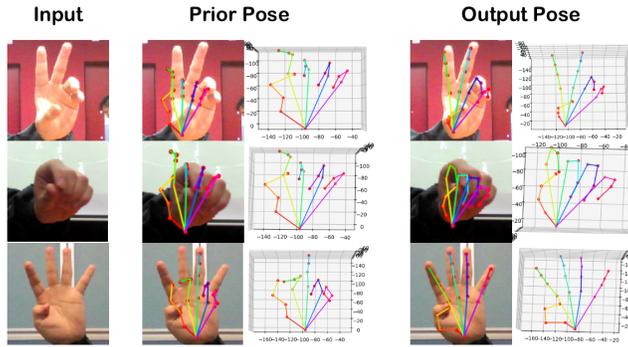| Input | Prior Pose | Output Pose |
|-------|-----------|-------------|



**Fig. 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.

crucial role of hand structure in pose estimation, but may resort to an additional synthetic dataset. Also, unlike bodies and faces that have obvious local characteristics (*e.g.*, eyes on a face), hands exhibit almost uniform appearance. Consequently, estimated hand poses from existing methods are sometimes distorted and unnatural.

To fully exploit the structure of hands, we propose to represent the irregular topology of 3D hand poses naturally on graphs, and learn the graph representation regularized by a prior pose from the monocular image input under a conditional generative adversarial learning framework, aiming to capture the structural dependencies among hand joints. Moreover, while most existing works [4,12,5] deploy 3D Euclidean distance between joints as the loss function for 3D annotation, we propose two *bone loss functions* that constrain the length and orientation of each bone connected by adjacent joints so as to preserve hand structure explicitly. Besides, unlike some recent works [12,5,18], we estimate 3D hand poses *without* resorting to ground truth meshes or depth maps, which is more suitable for datasets in the wild.

Specifically, given an input monocular image, our framework consists of a hand pose generator and a conditional discriminator. The generator is composed of a MANO hand model module [26] that provides an initial pose estimation as prior pose and a deformation learning module regularized by the prior pose. In particular, taking the prior pose and image features as input, the deformation learning module learns the deformation in the prior pose to further refine the hand structure, by our designed residual graph convolution that leverages on the recently proposed ResGCN [19]. Further, we design a conditional multi-source discriminator that employs hand poses, hand bones computed from poses as well as the input image to distinguish the predicted 3D hand pose from the ground-truth, leading to anthropomorphically valid hand pose. Experimental results demonstrate that our model achieves significant improvements over state-of-the-art approaches on five standard benchmarks.

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 information.
- 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 two categories: 1) 3D hand pose estimation from depth images; 2) 3D hand pose estimation from a monocular RGB image.

### 2.1    Estimation from Depth Images

Depth images contain rich 3D information for hand pose estimation [28], which has shown promising accuracy [32]. There is a rich literature on 3D hand pose estimation with depth images as input [10,11,9,6,7,16,21]. Among them, some earlier works such as [7,16] are based on a deformable hand model with an iterative optimization training approach. Due to the effectiveness of deep learning, some recent works like [21] leverage CNN to learn the shape and pose parameters for a proposed model (LBS hand model).

### 2.2    Estimation from a Monocular Image

Compared with the aforementioned two categories, a monocular RGB image is more accessible. Early works [2] propose complex model-fitting approaches, which are based on dynamics and multiple hypotheses and depend on restricted requirements. These model-fitting approaches have proposed many hand models, based on assembled geometric primitives [23] or sphere meshes [29], *etc.* Our work deploys the MANO hand model [26] as our prior, which models both hand shape and pose as well as generates meshes. Nevertheless, these sophisticated approaches suffer from low estimation accuracy.

   With the advance of deep learning, many recent works estimate 3D hand pose from a monocular RGB image using neural networks [12,4,3,5,18]. Among them, some recent works [18,12] directly reconstruct the 3D hand mesh and then generate the 3D hand pose through a pose regressor. Kulon *et al.* [18] reconstruct the hand pose based on an auto-encoder, which employs an encoder to extract the latent code and feeds the latent code into the decoder to reconstruct hand mesh. Ge *et al.* [12] propose to estimate vertices of 3D meshes from GCNs [17] in order to learn nonlinear variations in hand shape. The latent feature of the input RGB image is extracted via several networks and then fed into a GCN to directly infer the 3D coordinates of mesh vertices. However, since the accuracy of the
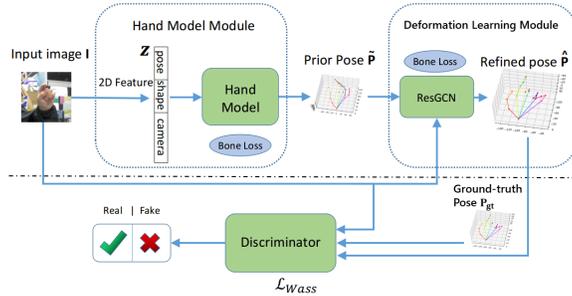
**Fig. 2.** Architecture of the proposed regularized graph representation learning under a conditional adversarial learning framework for 3D hand pose estimation.

output hand mesh is critical for both methods, they need an extra dataset which provides ground truth hand meshes as supervision. Also, the upsampling layer used in [12] 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 parametric hand model as regularization for graph representation learning over hand poses rather than directly reconstructing hand poses from latent features. Besides, our method does not require any additional supervision such as mesh supervision [12,18] or depth image supervision [12,5]. Hence, our method is more suitable for datasets in the wild. Further, we introduce conditional adversarial training for 3D hand pose estimation, which enables learning a real distribution of 3D hand poses.

## 3    Methodology

### 3.1    Overview of the Proposed Approach

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 $\mathbb{G}$ and a conditional discriminator $\mathbb{D}$, as illustrated in Fig. 2.

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 ground-truth 3D poses from the predicted ones.

### 3.2    The Proposed Hand Pose Generator $\mathbb{G}$

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

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

where $\mathbf{P}$ 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
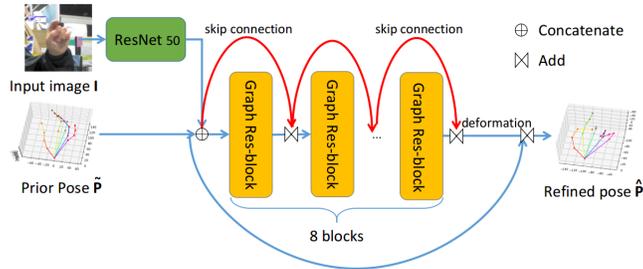
**Fig. 3.** Architecture of the deformation learning module in our generator.

**P**. $f(\mathbf{I}, \mathbf{P}_{\mathrm{gt}}|\mathbf{P})$ denotes the likelihood function, which is the probability of obtaining the observed image $\mathbf{I}$ and ground truth hand pose $\mathbf{P}_{\mathrm{gt}}$ given the estimated hand pose $\mathbf{P}$.

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

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

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

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

$d_1(\cdot)$ 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 $\mathbf{P}$, 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.

**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. Specifically, we crop and resize the input image to a salient region of the hand, which is fed into the ResNet-50 network [13] to extract features for the construction of the latent code $\mathbf{z}$, *i.e.*, parameters of the hand model. Then, we employ a modified MANO hand model [26], which is based on the SMPL model [20] for human bodies.

**The Deformation Learning Module** This module aims at accurate graph representation learning for hand pose estimation, which is conditional on the
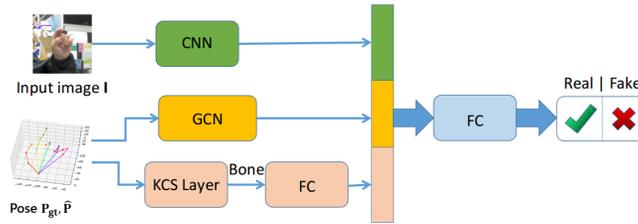
**Fig. 4.** Architecture of the conditional discriminator.

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.

We first construct an unweighted graph over $\tilde{\mathbf{P}}$, where the irregularly sampled key points (*i.e.*, joints) on the hand are projected onto nodes. The graph signal on each node is the concatenation of the global feature vector of the input image and the 3-dimensional coordinate vector of each joint in the input prior pose. Nodes are connected if they represent adjacent key points of the hand, where the adjacency relations follow the human hand structure as presented in Fig. 5, leading to an adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$.

Based on the graph representation $\mathbf{A}$, the finally refined pose is

$$\hat{\mathbf{P}} = \tilde{\mathbf{P}} + \text{GCN}(\tilde{\mathbf{P}} \oplus \mathbf{F}, \mathbf{A}), \tag{4}$$

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. $\text{GCN}(\tilde{\mathbf{P}} \oplus \mathbf{F}, \mathbf{A})$ represents the learned deformation between the prior $\tilde{\mathbf{P}}$ and the ground truth. The sum of the prior pose $\tilde{\mathbf{P}}$ and its deformation thus leads to the refined hand pose.

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

$$\mathbf{X}^{l+1} = N\left(g(N(g(\mathbf{X}^l, \mathbf{A})), \mathbf{A}))\right) + \text{skip}(\mathbf{X}^l), \tag{5}$$

where $g(\cdot)$ represents a single GCN layer as in [17], $N(\cdot)$ represents a single normalization layer, and $\text{skip}(\cdot)$ denotes the skip connection which is a GCN layer to match the dimension of the two terms in (5). We then stack several layers of Graph Res-blocks to learn the deformation of the prior pose, as demonstrated in Fig. 3.

### 3.3 The Proposed Conditional Discriminator $\mathbb{D}$

A simple architecture of a discriminator is a fully-connected (FC) network with the hand pose as input, which however has two shortcomings: 1) the relation between the RGB image and inferred hand pose is neglected; 2) structural properties of the hand pose are not taken into account explicitly. Instead, inspired
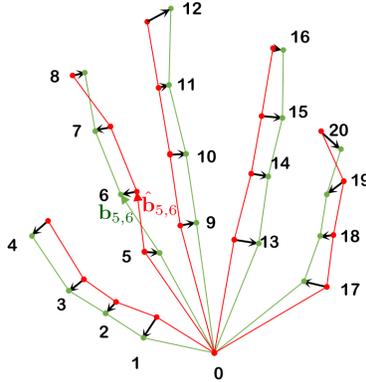
**Fig. 5. Illustration of the residual between the ground truth hand pose (marked in green) and the predicted one (marked in red).** Each hand pose has 21 key joints. We denote a bone vector connecting two key joints $i$ and $j$ by $\mathbf{b}_{i,j}$, such as $\mathbf{b}_{5,6}$ in the figure.

by the multi-source architecture in [31], we design a conditional multi-source discriminator with three inputs to address the aforementioned issues. As illustrated in Fig. 4, the inputs include: 1) features of the input monocular image; 2) features of the refined hand pose $\hat{\mathbf{P}}$ or the ground truth pose $\mathbf{P}_{\text{gt}}$; 3) features of bones via the KCS layer as in [30], which computes the bone matrix from $\hat{\mathbf{P}}$ or $\mathbf{P}_{\text{gt}}$ via a simple matrix multiplication. The bone features contain prominent structural information such as the length and direction of bones, thus characterizing the hand structure accurately.

The loss function of the conditional discriminator follows the definition of the Wasserstein loss [1] conditioned on the input image $\mathbf{I}$:

$$\mathcal{L}_{\text{Wass}} = -\mathbb{E}_{\mathbf{P}_{\text{gt}} \sim p_{data}(\mathbf{P}_{\text{gt}})} \mathbb{D}(\mathbf{P}_{\text{gt}}|\mathbf{I}) + \mathbb{E}_{\hat{\mathbf{P}} \sim p(\hat{\mathbf{P}})} \mathbb{D}(\hat{\mathbf{P}}|\mathbf{I}), \tag{6}$$

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

### 3.4   The Proposed Bone-Constrained Loss Functions

As presented in (3), we have two types of loss functions for the MAP estimation of hand pose. We employ the commonly adopted Euclidean distance in the coordinates of joints of 3D hand pose $\mathcal{L}_{\text{pose}}$ [12] as well as two proposed bone-constrained metrics as $d_1(\cdot)$ to measure the distortion of the estimated 3D hand pose compared to the ground truth, and apply the commonly used Euclidean distance in the coordinates of joints of projected 2D hand pose $\mathcal{L}_{\text{proj}}$ [12] as $d_2(\cdot)$ to measure the distance between the estimation and the 2D image,

$$\mathcal{L}_{\text{pose}} = \sum_i ||\mathbf{j}_i - \hat{\mathbf{j}}_i||_2, \mathcal{L}_{\text{proj}} = \sum_i ||\mathbf{j}'_i - \hat{\mathbf{j}}'_i||_2, \tag{7}$$

where $\mathbf{j}_i \in \mathbb{R}^{3\times1}, \mathbf{j}'_i \in \mathbb{R}^{2\times1}$ are 3D and 2D coordinates of joint $i$ respectively.

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

As illustrated in Fig. 5, we first define a bone vector $\mathbf{b}_{i,j} \in \mathbb{R}^{3\times1}$ between hand joint $i$ and $j$ as

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

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

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

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

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

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

Besides, as we adopt the framework of adversarial learning, we also introduce the Wasserstein loss $\mathcal{L}_{\text{Wass}}$ in (6) into the loss function for adversarial training. Hence, the overall loss function $\mathcal{L}$ is

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

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

## 4    Experimental Results

### 4.1    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.

**Stage II.** We train the generator $\mathbb{G}$ end-to-end without the discriminator $\mathbb{D}$. In $\mathbb{G}$, the hand model module is initialized with the trained model in the first stage and the deformation learning module is randomly initialized. $\mathbb{G}$ is then trained with 100 epochs using the Adam optimizer with learning rate 0.0001.

**Stage III.** We adopt the framework of SNGAN [22] for the conditional adversarial training, and train our model end-to-end. $\mathbb{G}$ and $\mathbb{D}$ are trained with 100 epochs using the Adam optimizer with learning rate 0.0001.

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

|  | STB | RHD | MPII+ZNSL(px) | Dexter+Object | EgoDexter |
|---|---|---|---|---|---|
| [12] | 6.37 | 15.33 | - | - | - |
| [4] | 9.76 | - | 18.95 | 25.53 | 45.33 |
| [27] | 8.56 | 19.73 | - | 40.20 | 56.92 |
| [34] | - | - | 59.40 | 34.75 | 52.77 |
| Ours | **3.97** | **12.40** | **9.87** | **16.12** | **34.98** |

**Table 1.** 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.

| Stage | hand model | deformation | discriminator | STB | RHD | EGODEXTER |
|---|---|---|---|---|---|---|
| I | ✓ |  |  | 24.15 | 83.37 | 52.32 |
| II | ✓ | ✓ |  | 5.12 | 15.84 | 43.26 |
| III | ✓ | ✓ | ✓ | **3.97** | **12.40** | **34.98** |

**Table 2.** The performance of different stages in our model on three datasets (measured in 3D Euclidean distance (mm)).

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

**Table 3.** 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).

### 4.2 Experimental Results

We compare our method with competitive 3D hand pose estimation approaches on the five datasets. We list the results in 3D Euclidean distance for comparison with the state-of-the-arts in Tab. 1. Compared to these works which directly reconstruct the 3D hand pose [12,4,5], our method performs much better mainly due to the proposed regularized graph representation learning and conditional adversarial learning. We show the qualitative results and PCK results in the supplementary material.

### 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 training stages in average 3D Euclidean distance, as listed in Tab. 2. The performance of **Stage II** significantly outperforms **Stage I**, which demonstrates that the proposed deformation learning module plays the most critical role in our model. The adversarial training scheme (**Stage III**) further improves the result, by learning a real distribution of the 3D hand pose.

| Model | Deformation Learning | Multi-source | Single-source | STB | RHD | EGODEXTER |
|-------|---------------------|--------------|---------------|------|-------|-----------|
| 1 | ✓ | ✓ | | **3.97** | **12.40** | **34.98** |
| 2 | ✓ | | ✓ | 4.54 | 15.10 | 37.46 |

**Table 4.** Ablation studies on the discriminator (3D Euclidean distance (mm)).

| Model | $\mathcal{L}_{\mathrm{pose}} + \mathcal{L}_{\mathrm{proj}}$ | $\mathcal{L}_{\mathrm{len}}$ | $\mathcal{L}_{\mathrm{dir}}$ | STB | | | RHD | | |
|-------|:---:|:---:|:---:|---------|----------|-----------|---------|----------|-----------|
| | | | | 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 | ✓ | ✓ | | 30.32 | 8.00 | 5.02 | 95.19 | 22.96 | 14.76 |
| 3 | ✓ | | ✓ | 27.65 | 6.91 | 5.00 | 89.76 | 21.63 | 14.01 |
| 4 | ✓ | ✓ | ✓ | 24.15 | 5.12 | **3.97** | 83.37 | 15.84 | **12.40** |

**Table 5.** Ablation studies on the proposed bone-constrained loss functions at three stages.

**On the deformation learning module.** We compare the deformation learning module with a simple fully-connected deformation learning module (FC Deformation Module) to refine the prior pose. We train the deformation learning modules in different experimental settings: 1) without our discriminator, *i.e.*, without adversarial learning; and 2) with our discriminator. As presented in Tab. 3, the GCN deformation learning module leads to significant gain over the simple FC deformation module on both datasets in different settings, thus validating the superiority of the proposed deformation learning module.

**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. 4, 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.

**On loss functions.** We also evaluate the proposed bone-constrained loss functions $\mathcal{L}_{\mathrm{len}}$ and $\mathcal{L}_{\mathrm{dir}}$ separately. We train the network with different combinations of loss functions on the STB and RHD datasets in three stages respectively. As reported in Tab. 5, the network trained with our proposed bone-constrained loss functions performs better in all the three stages on both datasets. We also notice that $\mathcal{L}_{dir}$ plays a more significant role compared to $\mathcal{L}_{\mathrm{len}}$. This gives credits to the constraint on the orientation of bones that explicitly takes structural properties of hands into consideration.

## 5   Conclusion

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

## References

1. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein generative adversarial networks. In: Precup, D., Teh, Y.W. (eds.) Proceedings of the 34th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 70, pp. 214–223. PMLR, International Convention Centre, Sydney, Australia (06–11 Aug 2017)
2. Athitsos, V., Sclaroff, S.: Estimating 3d hand pose from a cluttered image. In: IEEE Computer Society Conference on Computer Vision & Pattern Recognition (2003)
3. Baek, S., Kim, K.I., Kim, T.K.: Pushing the envelope for rgb-based dense 3d hand pose estimation via neural rendering. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
4. Boukhayma, A., Bem, R.d., Torr, P.H.: 3d hand shape and pose from images in the wild. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
5. Cai, Y., Ge, L., Cai, J., Yuan, J.: Weakly-supervised 3d hand pose estimation from monocular rgb images. In: The European Conference on Computer Vision (ECCV) (September 2018)
6. Choi, C.: Deephand: Robust hand pose estimation by completing a matrix imputed with deep features. In: Computer Vision & Pattern Recognition (2016)
7. De, L.G.M., Fleet, D.J., Paragios, N.: Model-based 3d hand pose estimation from monocular video. IEEE Trans Pattern Anal Mach Intell **33**(9), 1793–1805 (2011)
8. Doosti, B., Naha, S., Mirbagheri, M., Crandall, D.J.: Hope-net: A graph-based model for hand-object pose estimation. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020)
9. Fitzgibbon, A.: Accurate, robust, and flexible real-time hand tracking. Inproceedings pp. 3633–3642 (2015)
10. Ge, L., Cai, Y., Weng, J., Yuan, J.: Hand pointnet: 3d hand pose estimation using point sets. pp. 8417–8426 (06 2018). https://doi.org/10.1109/CVPR.2018.00878
11. Ge, L., Liang, H., Yuan, J., Thalmann, D.: Robust 3d hand pose estimation in single depth images: From single-view cnn to multi-view cnns. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016)
12. Ge, L., Ren, Z., Li, Y., Xue, Z., Wang, Y., Cai, J., Yuan, J.: 3d hand shape and pose estimation from a single rgb image. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016)
14. Hui, L., Yuan, J., Lee, J., Ge, L., Thalmann, D.: Hough forest with optimized leaves for global hand pose estimation with arbitrary postures. IEEE Transactions on Cybernetics **PP**(99), 1–15 (2017)
15. Hürst, W., van Wezel, C.: Gesture-based interaction via finger tracking for mobile augmented reality. Multimedia Tools and Applications **62**, 233–258 (2011)
16. Khamis, S., Taylor, J., Shotton, J., Keskin, C., Izadi, S., Fitzgibbon, A.: Learning an efficient model of hand shape variation from depth images. In: IEEE Conference on Computer Vision & Pattern Recognition (2015)
17. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: International Conference on Learning Representations (ICLR) (2017)
18. Kulon, D., Wang, H., Güler, R.A., Bronstein, M.M., Zafeiriou, S.: Single image 3d hand reconstruction with mesh convolutions. In: BMVC (September 2019)

19. Li, G., Muller, M., Thabet, A., Ghanem, B.: Deepgcns: Can gcns go as deep as cnns? In: The IEEE International Conference on Computer Vision (ICCV) (October 2019)
20. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: Smpl: A skinned multi-person linear model. ACM Trans. Graph. **34**(6), 248:1–248:16 (Oct 2015). https://doi.org/10.1145/2816795.2818013, `http://doi.acm.org/10.1145/2816795.2818013`
21. Malik, J., Elhayek, A., Nunnari, F., Varanasi, K., Stricker, D.: Deephps: End-to-end estimation of 3d hand pose and shape by learning from synthetic depth. In: 2018 International Conference on 3D Vision (3DV) (2018)
22. Miyato, T., Kataoka, T., Koyama, M., Yoshida, Y.: Spectral normalization for generative adversarial networks. In: International Conference on Learning Representations (2018), `https://openreview.net/forum?id=B1QRgziT-`
23. Oikonomidis, I., Kyriazis, N., Argyros, A.: Efficient model-based 3d tracking of hand articulations using kinect. vol. 1 (01 2011). https://doi.org/10.5244/C.25.101
24. Panteleris, P., Argyros, A.A.: Back to RGB: 3d tracking of hands and hand-object interactions based on short-baseline stereo. CoRR **abs/1705.05301** (2017), `http://arxiv.org/abs/1705.05301`
25. Piumsomboon, T., Clark, A., Billinghurst, M., Cockburn, A.: User-defined gestures for augmented reality. In: Kotzé, P., Marsden, G., Lindgaard, G., Wesson, J., Winckler, M. (eds.) Human-Computer Interaction – INTERACT 2013. pp. 282–299. Springer Berlin Heidelberg, Berlin, Heidelberg (2013)
26. Romero, J., Tzionas, D., Black, M.J.: Embodied hands: Modeling and capturing hands and bodies together. ACM Trans. Graph. **36**(6), 245:1–245:17 (Nov 2017). https://doi.org/10.1145/3130800.3130883, `http://doi.acm.org/10.1145/3130800.3130883`
27. Spurr, A., Song, J., Park, S., Hilliges, O.: Cross-modal deep variational hand pose estimation. CoRR **abs/1803.11404** (2018), `http://arxiv.org/abs/1803.11404`
28. Tang, D., Yu, T.H., Kim, T.K.: Real-time articulated hand pose estimation using semi-supervised transductive regression forests. In: IEEE International Conference on Computer Vision (2013)
29. Tkach, A., Pauly, M., Tagliasacchi, A.: Sphere-meshes for real-time hand modeling and tracking. ACM Trans. Graph. **35**(6), 222:1–222:11 (Nov 2016). https://doi.org/10.1145/2980179.2980226, `http://doi.acm.org/10.1145/2980179.2980226`
30. Wandt, B., Ackermann, H., Rosenhahn, B.: A kinematic chain space for monocular motion capture (02 2017)
31. Yang, W., Ouyang, W., Wang, X., Ren, J., Li, H., Wang, X.: 3d human pose estimation in the wild by adversarial learning. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
32. Yuan, S., Garcia-Hernando, G., Stenger, B., Moon, G., Chang, J., Lee, K., Molchanov, P., Kautz, J., Honari, S., Ge, L., Yuan, J., Chen, X., Wang, G., Yang, F., Akiyama, K., Wu, Y., Wan, Q., Madadi, M., Escalera, S., Kim, T.K.: Depth-based 3d hand pose estimation: From current achievements to future goals. pp. 2636–2645 (06 2018). https://doi.org/10.1109/CVPR.2018.00279
33. Zhang, J., Jiao, J., Chen, M., Qu, L., Xu, X., Yang, Q.: 3d hand pose tracking and estimation using stereo matching (10 2016)
34. Zimmermann, C., Brox, T.: Learning to estimate 3d hand pose from single rgb images. In: The IEEE International Conference on Computer Vision (ICCV) (Oct 2017)