

MobRecon: Mobile-Friendly Hand Mesh Reconstruction from Monocular Image

Xingyu Chen^{1*} Yufeng Liu³ Yajiao Dong¹ Xiong Zhang² Chongyang Ma¹
Yanmin Xiong¹ Yuan Zhang¹ Xiaoyan Guo¹

¹Y-tech, Kuaishou Technology ²YY Live, Baidu Inc.

³SEU-ALLEN Joint Center, Institute for Brain and Intelligence, Southeast University, China.

Abstract

In this work, we propose a framework for single-view hand mesh reconstruction, which can simultaneously achieve high reconstruction accuracy, fast inference speed, and temporal coherence. Specifically, for 2D encoding, we propose lightweight yet effective stacked structures. Regarding 3D decoding, we provide an efficient graph operator, namely depth-separable spiral convolution. Moreover, we present a novel feature lifting module for bridging the gap between 2D and 3D representations. This module begins with a map-based position regression (MapReg) block to integrate the merits of both heatmap encoding and position regression paradigms for improved 2D accuracy and temporal coherence. Furthermore, MapReg is followed by pose pooling and pose-to-vertex lifting approaches, which transform 2D pose encodings to semantic features of 3D vertices. Overall, our hand reconstruction framework, called MobRecon, comprises affordable computational costs and miniature model size, which reaches a high inference speed of 83FPS on Apple A14 CPU. Extensive experiments on popular datasets such as FreiHAND, RHD, and HO3Dv2 demonstrate that our MobRecon achieves superior performance on reconstruction accuracy and temporal coherence. Our code is publicly available at <https://github.com/SeanChenxy/HandMesh>.

1. Introduction

Single-view hand mesh reconstruction has been extensively investigated for years due to its wide range of applications in AR/VR [28, 73], behavior understanding [40, 63], etc. Tremendous research efforts have been made towards this task, including [18, 90, 46, 85], to name a few.

The primary focus of typical existing methods is the reconstruction accuracy [50, 51], while real-world applications additionally necessitate inference efficiency and temporal consistency. In particular, 3D hand information

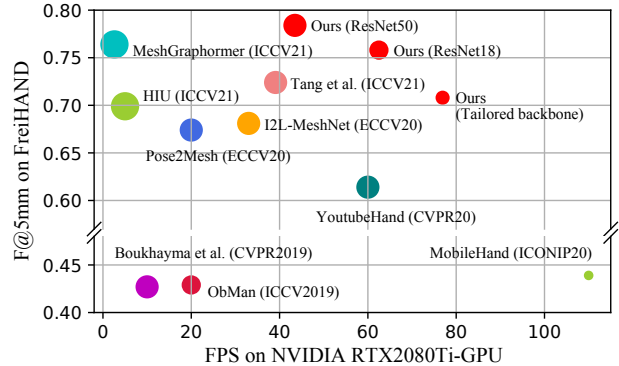


Figure 1. Accuracy vs. inference speed. The marker size is related to the model size. Besides, our tailored method can run at a fast speed on mobile CPUs.

is a vital component in mobile applications [73], where the devices comprise limited memory and computational budgets. Thereby, this work aims to explore 3D hand reconstruction for mobile platforms.

A typical pipeline for single-view hand reconstruction includes three phases: 2D encoding, 2D-to-3D mapping, and 3D decoding. In 2D encoding, existing approaches [46, 11, 50, 51] usually adopt computationally intensive networks [30, 75] to handle this highly non-linear task, which are hard to deploy on mobile devices. Instead, if naively leveraging a mature mobile network (e.g., [32]) which is not tailored for our target task, the reconstruction accuracy dramatically degrades [22]. Hence, our motivation is to develop a lightweight 2D encoding structure tailored to balance the inference efficiency and accuracy. Besides, the efficiency of 2D-to-3D mapping and 3D decoding remains relatively unexplored. Thus, we intend to explore a lightweight yet effective lifting method to tackle the 2D-to-3D mapping problem and design an efficient graph operator for processing of 3D mesh data.

Although as crucial as accuracy in real-world applications, temporal coherence is usually neglected in the task of 3D hand reconstruction. Previous methods [13, 41, 45, 54] adopt sequential models to incorporate both past and future

*Corresponding author, chenxingyu@kuaishou.com

semantic information for stable predictions. Since this methodology is offline or computationally expensive, these approaches are difficult to use for mobile applications. Hence, we are inspired to explore the temporal coherence with a non-sequential method.

In this work, we propose **Mobile Mesh Reconstruction** (MobRecon) for 3D hand to simultaneously explore superior accuracy, efficiency, and temporal coherence. For 2D encoding, we leverage the spirit of the hourglass network [62] to design efficient stacked encoding structures. As for 3D decoding, we propose a depth-separable spiral convolution (DSConv), which is a novel graph operator based on spiral sampling [49]. The DSConv is inspired by depth-wise separable convolution [33], leading to efficient handling of graph-structured mesh data. Regarding the 2D-to-3D mapping, we propose a feature lifting module with map-based position regression (MapReg), pose pooling, and pose-to-vertex lifting (PVL) approaches. In this module, we first investigate the pros and cons of 2D pose representations based on heatmap or position regression, and then propose a hybrid method MapReg to simultaneously improve 2D pose accuracy and temporal consistency. Furthermore, the PVL transforms 2D pose encodings to 3D vertex features based on a learnable lifting matrix, resulting in enhanced 3D accuracy and temporal consistency. Compared to traditional approaches based on fully connected operation in a latent space [46, 18, 11], our feature lifting module also significantly reduces the model size. In addition, we build a synthetic dataset with uniformly distributed hand poses and viewpoints. Referring to Figure 1, we achieve better performance in terms of accuracy, speed, and model size.

Our main contributions are summarized as follows:

- We propose MobRecon as a mobile-friendly pipeline for hand mesh reconstruction, which only involves 123M multiply-add operations (Mult-Adds) and 5M parameters and can run up to 83FPS on Apple A14 CPU.
- We present lightweight stacked structures and DSConv for efficient 2D encoding and 3D decoding.
- We propose a novel feature lifting module with MapReg, pose pooling, and PVL methods to bridge the 2D and 3D representations.
- We demonstrate that our method achieves superior performance in terms of model efficiency, reconstruction accuracy, and temporal coherence via comprehensive evaluations and comparisons with state-of-the-art approaches.

2. Related Work

Hand mesh estimation. Popular hand mesh estimation methods can be divided into five types, whose core ideas are based on the parametric model, voxel representation, implicit function, UV map, and vertex position, respectively.

Model-based approaches [86, 90, 80, 29, 85, 4, 1, 93, 12, 88, 6, 2, 39, 52, 82, 83] typically use MANO [65] as the parametric model, which factorizes a hand mesh into coefficients of shape and pose. This pipeline, however, is not suitable for usage with a lightweight network because the coefficient estimation is a highly abstract problem that ignores spatial correlations [22].

Voxel-based approaches [37, 57, 79, 59] describe the 3D data in a 2.5D manner. Moon *et al.* [57] proposed I2L-MeshNet, which divided the voxel space into three lixel spaces and used 1D heatmaps to reduce memory consumption. Despite this optimization, I2L-MeshNet still requires massive memory occupation to process the lixel-style 2.5D heatmaps. Hence, the voxel-based approach is not friendly for memory-constrained mobile devices.

The implicit function [55] has merits of continuity and high resolution, which is recently used for digitizing articulated human [56, 38, 3, 66, 36, 64, 42, 43]. However, the implicit methods usually need to compute thousands of 3D points, lacking efficiency in mobile settings.

Chen *et al.* [10] treated hand mesh reconstruction as an image-to-image translation task, employing UV map to connect 2D and 3D spaces. This pipeline could be improved by incorporating geometry correlations.

Vertex-based methods [18, 46, 11, 50, 51] predict 3D vertex coordinates directly, which usually follow a procedure of 2D encoding, 2D-to-3D mapping, and 3D decoding. For example, Kulon *et al.* [46] designed an encoder-decoder based on ResNet [30], global pooling, and spiral convolution (SpiralConv) [49] to obtain 3D vertex coordinates. We re-construct the vertex-based pipeline with efficient modules, *i.e.*, lightweight stacked structures for 2D encoding, feature lifting module for 2D-to-3D mapping, and DSConv for 3D decoding. As a result, we achieve high reconstruction accuracy and across-time consistency.

Lightweight networks. For a timely fashion on a computationally limited platform, lightweight networks have been studied for years such as [33, 67, 32, 27]. We leverage popular efficient ideas to design stacked networks for Euclidean 2D images and a graph network for non-Euclidean 3D meshes. More specifically, we propose a feature lifting module to deal with the problem of 2D-to-3D mapping efficiently. As a highly related work, MobileHand [22] was able to run at 75FPS on a mobile CPU. In contrast, our MobRecon achieves more powerful performance on accuracy and inference speed (as shown in Table 5).

Temporally coherent human reconstruction. There has been limited research on the temporal coherence of human/hand mesh reconstruction, yet it is as crucial as reconstruction accuracy in real-world applications. Previous research [13, 41, 45, 54] has concentrated on temporal performance with temporal approaches. Kocabas *et al.* [45] used a bi-directional gated recurrent unit [14] to fuse

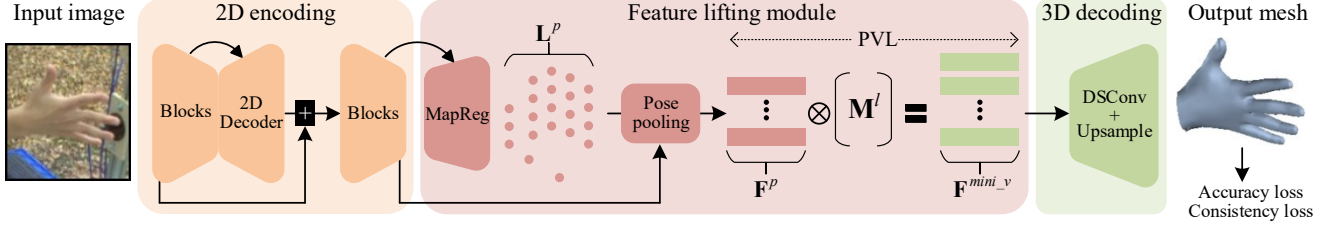


Figure 2. Overview of our MobRecon framework.

across-time features so that SMPL [53] parameters could be regressed with temporal cues. This sequential manner could raise computational costs and even requires future information [13]. In contrast, we design a feature lifting module with MapReg to enhance the temporal coherence for a non-sequential single-view method.

Pixel-aligned representation. Convolutional features are made up of dense and regularly structured 2D cues, but sparse and unordered points represent 3D data. To better extract image features to describe 3D information, recent works usually adopt pixel-aligned representations [66, 58, 23, 83, 26, 89]. Inspired by them, we use the idea of pixel alignment for feature lifting and design PVL to transform pose-aligned encodings to vertex features.

It is noteworthy that the above-mentioned literature used heatmap [23, 26, 83] or position [66, 89] as the 2D representations. Li *et al.* [48] analyzed the heatmap- and position-based human pose in terms of accuracy and proposed RLE to achieve high-accuracy regression. We consider these two representations from the view of temporal coherence and propose MapReg to integrate the merits of heatmap- and position-based 2D representations. Compared to [48], we have different insights and study perspectives, so RLE and our MapReg could complement each other.

3. Our Method

With a single-view image as the input, we aim to infer 3D hand mesh with predicted vertices $\mathbf{V} = \{\mathbf{v}_i\}_{i=1}^V$ and pre-defined faces $\mathbf{C} = \{\mathbf{c}_i\}_{i=1}^C$. Figure 2 illustrates the overall architecture of MobRecon that includes three phases. For 2D encoding, we leverage convolutional networks to extract image features. In the feature lifting module, 2D pose is delineated with MapReg, followed by pose pooling to retrieve pose-aligned features. Then, vertex features are obtained with PVL. For 3D decoding, we develop an efficient graph network to predict vertex coordinates.

3.1. Stacked Encoding Network

Inspired by the hourglass network [62], we develop a stacked encoding network to obtain gradually refined encoding features. As shown in Figure 3, the stacked network consists of two groups of cascaded encoding blocks, the first of which is followed by an upsampling module for feature

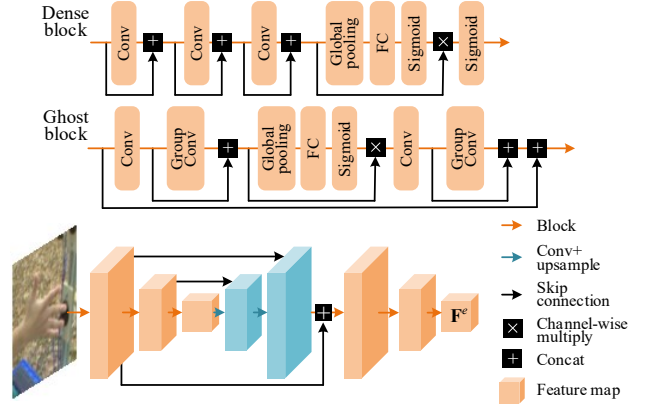


Figure 3. Our tailored stacked encoding structure.

fusion. With a single-view image as the input, the encoding feature $\mathbf{F}^e \in \mathbb{R}^{C^e \times H^e \times W^e}$ is generated, where C, H, W denote tensor channel size, height, and width.

We design two block alternatives. According to DenseNet [35] and SENet [34], we present a dense block to form DenseStack (see Figure 3). To further reduce model size, we leverage the ghost operation [27] to develop GhostStack, where cheap operations can produce ghost features based on primary features. With 128×128 input resolution, the DenseStack involves 373.0M Mult-Adds and 6.6M parameters while the GhostStack consists of 96.2M Mult-Adds and 5.0M parameters. In contrast, a stacked network with ResNet18 has 2391.3M Mult-Adds and 25.2M parameters, prohibiting it from mobile applications.

3.2. Feature Lifting Module

Lifting means 2D-to-3D mapping [80]. For feature lifting, two problems should be concerned: (1) how to collect 2D features and (2) how to map them to 3D domain. To this end, previous methods [46, 18, 11] tend to embed \mathbf{F}^e as a latent vector via the global average pooling operation. Then, the latent vector is mapped to 3D domain with a fully connected layer (FC), and vertex features are obtained with vector re-arrangement. This manner brings increased model size due to the large dimension of FC, *i.e.*, this layer contains 3.2M parameters when $C^e=256$.

Recent researches report pixel-aligned feature extraction based on 2D landmarks and pixel-aligned feature pooling [66, 23, 26, 83, 89]. Heatmap \mathbf{H}^p is usually employed to

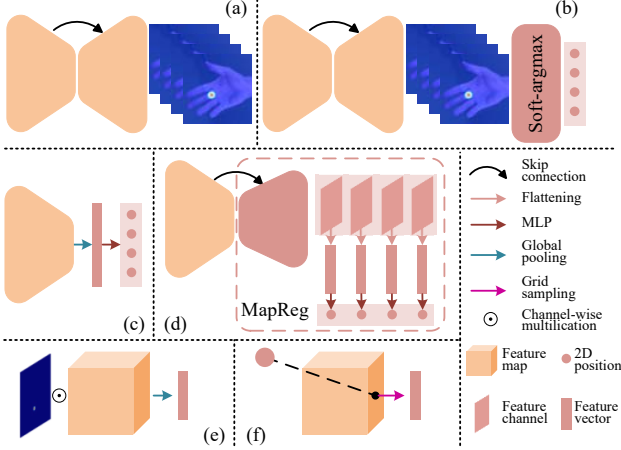


Figure 4. Comparisons of 2D representations and pose pooling methods. (a) heatmap, (b) heatmap + soft-argmax, (c) regression-based position, (d) MapReg-based position, (e) joint-wise pooling with a heatmap, (f) grid sampling with a 2D position. For better visualization, only 4 landmarks are illustrated.

encode 2D landmarks [23, 18, 11, 85], which derives more accurate landmarks compared with direct regression of the 2D positions \mathbf{L}^p [48, 74, 7, 71, 76].

Map-based position regression. Figure 4 comprehensively investigates these 2D representations. It is seen that \mathbf{H}^p is a high-resolution representation (Figure 4(a)). As shown in Figure 4(b), soft-argmax [72] is a differentiable technique that decodes \mathbf{H}^p to the corresponding 2D position. In contrast, direct position regression (Figure 4(c)) is a low-resolution representation, which can yield \mathbf{L}^p without resorting to \mathbf{H}^p . It has been proven that the human pose estimation task requires both low- and high-resolution encodings [75]. Hence, the skip connection that fuses features is the primary cause behind the superior accuracy of \mathbf{H}^p [62]. However, our critical insight is that the global feature (with resolution of 1×1) used to predict \mathbf{L}^p (see Figure 4(c)) shall induce better temporal coherence because of global semantics and receptive field. This global property can better describe articulated relation of hand pose. Conversely, \mathbf{H}^p is predicted with convolution and high-resolution features, where limited receptive field leads to the lack of inter-landmark constraints.

As shown in Figure 4(d), we propose a middle-resolution method MapReg to combine the advantages of heatmap- and position-based paradigms by (1) fusing low- and high-resolution features for accuracy and (2) using global receptive field for temporal coherence. To this end, we incorporate the skip connection in the position regression paradigm, producing a spatially small-size (e.g. 16×16) feature map. Each feature channel is then flattened into a vector, followed by a multi-layer perceptron (MLP) to generate a 2D position. In this manner, we obtain middle-resolution spatial complexity which is superior than heatmap since

only two $2 \times$ -upsampling operations are involved.

Pose pooling. After obtaining 2D representations, pixel-aligned features can be retrieved. We call this process pose pooling and capture pose-aligned features with N 2D joint landmarks. If heatmap \mathbf{H}^p is employed as the 2D representation, the pose pooling is conducted with joint-wise pooling [83] (Figure 4(e)), which can be given as

$$\mathbf{F}^p = [\Psi(\mathbf{F}^e \odot \text{interpolation}(\mathbf{H}_i^p))]_{i=1,2,\dots,N}, \quad (1)$$

where $[\cdot]$ denotes concatenation. First, the spatial size of \mathbf{H}^p is interpolated to $H^e \times W^e$, and then channel-wise multiplication is adopted between interpolated \mathbf{H}^p and \mathbf{F}^e . In this manner, features unrelated to joint landmarks are suppressed. For extracting joint-wise features, feature reduction Ψ is designed. In detail, Ψ denotes global max-pooling or spatial sum reduction, which produces a C^e -length feature vector. After concatenation, $\mathbf{F}^p \in \mathbb{R}^{N \times C^e}$ indicates the pose-aligned feature.

If we use \mathbf{L}^p instead of \mathbf{H}^p to describe 2D pose, pose pooling can be achieved with grid sampling [66] (Figure 4(f)) as follows,

$$\mathbf{F}^p = [\mathbf{F}^e(\mathbf{L}_i^p)]_{i=1,2,\dots,N}. \quad (2)$$

As a result, the convolutional encoding \mathbf{F}^e is transformed to the pose-aligned representation \mathbf{F}^p .

Pose-to-vertex lifting. Referring to Figure 2, we design a linear operator for feature mapping towards 3D space with a few learnable parameters, namely PVL. MANO-style hand mesh [65] comprises V vertices and N joints, where $V = 778, N = 21$. Because $V \gg N$, it is hard to transform \mathbf{F}^p to V vertex features. Instead, we downsample the template mesh 4-times by a factor of 2 [17], and obtain a minimal-size hand mesh with $V^{mini} = 49$ vertices. Then, we design a learnable lifting matrix $\mathbf{M}^l \in \mathbb{R}^{V^{mini} \times N}$ for 2D-to-3D feature mapping. Thereby, PVL is given as

$$\mathbf{F}^{mini-v} = \mathbf{M}^l \cdot \mathbf{F}^p, \quad (3)$$

where \mathbf{F}^{mini-v} denotes minimal-size vertex feature. The PVL reduces the computational complexity of feature mapping from $\mathcal{O}(V^{mini} C^e 2)$ [46, 11, 18] to $\mathcal{O}(N V^{mini} C^e)$.

3.3. Depth-Separable SpiralConv

As a graph operator, SpiralConv [49] is equivalent to the Euclidean convolution, which designs a spiral neighbor as

$$\begin{aligned} 0\text{-ring}(\mathbf{v}) &= \{\mathbf{v}\} \\ (k+1)\text{-ring}(\mathbf{v}) &= \mathbb{N}(k\text{-ring}(\mathbf{v})) \setminus k\text{-disk}(\mathbf{v}) \\ k\text{-disk}(\mathbf{v}) &= \cup_{i=0,\dots,k} i\text{-ring}(\mathbf{v}), \end{aligned} \quad (4)$$

where \mathbb{N} extracts the neighbourhood of a vertex \mathbf{v} . With $k\text{-disk}(\mathbf{v})$, SpiralConv formulates convolution as a sequential problem and leverages LSTM [31] for feature fusion:

$$\mathbf{f}_{\mathbf{v}}^{out} = \text{LSTM}(\mathbf{f}_{\mathbf{v}'}), \quad \mathbf{v}' \in k\text{-disk}(\mathbf{v}), \quad (5)$$

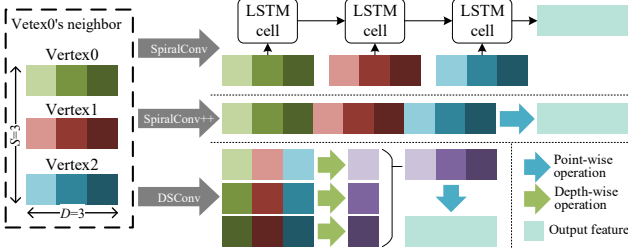


Figure 5. Comparison of SpiralConv, SpiralConv++, and DSConv. For better visualization, a case with $S = D = 3$ is shown.

where $\mathbf{f}_v \in \mathbb{R}^{1 \times D}$ denotes the feature vector at \mathbf{v} with dimension D . SpiralConv with LSTM could be potentially slow because of serial sequence processing.

By explicitly formulating the order of aggregating neighboring vertices, SpiralConv++ [20] presents an efficient version of SpiralConv. For high efficiency, SpiralConv++ only adopts a fixed-size spiral neighbor with S vertices and leverages FC to fuse these features:

$$\mathbf{f}_{\mathbf{v}'}^{\text{out}} = \mathbf{W} \cdot [\mathbf{f}_{\mathbf{v}'}] + \mathbf{b}, \quad \mathbf{v}' \in k\text{-disk}(\mathbf{v})_S, \quad (6)$$

where $[\cdot]$ denotes concatenation; $k\text{-disk}(\mathbf{v})_S$ contains the first S elements in $k\text{-disk}(\mathbf{v})$; \mathbf{W} and \mathbf{b} are learnable parameters. SpiralConv++ significantly increases model size because of the large dimension of FC.

As an efficient graph convolution, we propose DSConv with the spirit of depth-wise separable convolution [33]. For a vertex \mathbf{v} , $k\text{-disk}(\mathbf{v})_S$ is sampled following Equation 4. DSConv comprises a depth-wise operation and a point-wise operation, the former of which can be formulated as

$$\mathbf{f}_{\mathbf{v}'}^d = [\mathbf{W}_i^d \cdot [\mathbf{f}_{\mathbf{v}',i}]]_{i=1}^D, \quad \mathbf{v}' \in k\text{-disk}(\mathbf{v})_S. \quad (7)$$

Then, point-wise operation can be formulated as

$$\mathbf{f}_{\mathbf{v}}^{\text{out}} = \mathbf{W}^p \cdot \mathbf{f}_{\mathbf{v}}^d. \quad (8)$$

Figure 5 illustrates the difference among SpiralConv, SpiralConv++, and our DSConv. The SpiralConv++'s computational complexity obeys $\mathcal{O}(SD^2)$, whereas that of DSConv is $\mathcal{O}(SD + D^2)$. Hence, we essentially improve the efficiency with the separable structure.

The 3D decoder is built with four blocks, each of which involves upsampling, DSConv, and ReLU. In each block, vertex features are upsampled by a factor of 2 and then processed by DSConv. Finally, vertex coordinates \mathbf{V} are predicted by a DSConv.

3.4. Loss Functions

Accuracy loss. We use L_1 norm for the 3D mesh loss $\mathcal{L}_{\text{mesh}}$ and 2D pose loss $\mathcal{L}_{\text{pose2D}}$. Normal loss $\mathcal{L}_{\text{norm}}$ and edge length loss $\mathcal{L}_{\text{edge}}$ are adopted for mesh smoothness

according to [18]. Formally, we have

$$\begin{aligned} \mathcal{L}_{\text{mesh}} &= \|\mathbf{V} - \mathbf{V}^*\|_1, \quad \mathcal{L}_{\text{pose2D}} = \|\mathbf{L}^p - \mathbf{L}^{p,*}\|_1 \\ \mathcal{L}_{\text{norm}} &= \sum_{\mathbf{c} \in \mathbf{C}} \sum_{(i,j) \in \mathbf{c}} \left| \frac{\mathbf{V}_i - \mathbf{V}_j}{\|\mathbf{V}_i - \mathbf{V}_j\|_2} \cdot \mathbf{n}_{\mathbf{c}}^* \right| \\ \mathcal{L}_{\text{edge}} &= \sum_{\mathbf{c} \in \mathbf{C}} \sum_{(i,j) \in \mathbf{c}} \left| \|\mathbf{V}_i - \mathbf{V}_j\|_2 - \|\mathbf{V}_i^* - \mathbf{V}_j^*\|_2 \right|, \end{aligned} \quad (9)$$

where \mathbf{C} , \mathbf{V} are face and vertex sets of a mesh; $\mathbf{n}_{\mathbf{c}}^*$ indicates unit normal vector of face \mathbf{c} ; $*$ denotes the ground truth.

Consistency loss. Inspired by the self-supervision task [68, 91], we design the consistency supervision based on augmentation without the need of temporal data. That is, two views can be derived from an image sample based on 2D affine transformation and color jitter. We denote relative affine transformation between two views as $T_{1 \rightarrow 2}$, which contains relative rotation $R_{1 \rightarrow 2}$. Similar to [24, 79], we conduct consistency supervision in both 3D and 2D space:

$$\begin{aligned} \mathcal{L}_{\text{con3D}} &= \|R_{1 \rightarrow 2} \mathbf{V}_{\text{view1}} - \mathbf{V}_{\text{view2}}\|_1 \\ \mathcal{L}_{\text{con2D}} &= \|T_{1 \rightarrow 2} \mathbf{L}_{\text{view1}}^p - \mathbf{L}_{\text{view2}}^p\|_1. \end{aligned} \quad (10)$$

Although T contains the variations of rotation, shift, and scale in 2D space, only R affects \mathbf{V} because it is rooted by the wrist landmark.

Our overall loss function is $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{mesh}} + \mathcal{L}_{\text{pose2D}} + \mathcal{L}_{\text{norm}} + \mathcal{L}_{\text{edge}} + \mathcal{L}_{\text{con2D}} + \mathcal{L}_{\text{con3D}}$.

4. Experiments

4.1. Implementation Details

We use the Adam optimizer [44] to train the network with a mini-batch size of 32. All models in our experiments are trained for 38 epochs. The initial learning rate is 10^{-3} , which is divided by 10 at the 30th epoch. As for hyperparameters, the input resolution is 128×128 , and $S = 9, C^e = 256, D = \{256, 128, 64, 32\}$.

Existing datasets for 3D hand pose estimation (e.g. [93, 46, 25]) usually suffer from long-tailed distribution of hand pose and viewpoint. Hence, we develop a synthetic dataset with 1520 poses and 216 viewpoints, both of which are uniformly distributed in their respective spaces. Because of this superior property, it can serve as a good complement during training. We redirect the readers to supplementary material for details.

4.2. Evaluation Criterion

We conduct experiments on several commonly-used benchmarks as listed below.

FreiHAND [93] contains 130,240 training images and 3,960 evaluation samples. The annotations of the evaluation set are not available, so we submit our predictions to the official server for online evaluation.

2D encoding	Mult-Adds	#Param	Fine-tuning data	PJ↓
<i>From-scratch fine-tuning</i>				
ResNet18-Stack	2391.3M	25.2M	FreiHAND	8.86
MobileNet-Stack	100.2M	1.7M	FreiHAND	14.20
GhostStack	96.2M	5.0M	FreiHAND	13.09
DenseStack	373.0M	6.6M	FreiHAND	9.56
<i>Pre-training w/ ImageNet</i>				
ResNet18-Stack	2391.3M	25.2M	FreiHAND	8.21
MobileNet-Stack	100.2M	1.7M	FreiHAND	13.68
<i>Pre-training w/ ours</i>				
MobileNet-Stack	100.2M	1.7M	FreiHAND	12.45
			FreiHAND	10.05
GhostStack	96.2M	5.0M	FreiHAND+ours	8.89
			FreiHAND	7.77
DenseStack	373.0M	6.6M	FreiHAND+ours	7.55

Table 1. Ablation studies of 2D encoding methods and our complement data on FreiHAND. Mult-Adds and #Param are *w.r.t.* the 2D encoding network.

Rendered Hand Pose Dataset (RHD) [92] consists of 41,258 and 2,728 synthetic hand data for training and testing on hand pose estimation, respectively.

HO3Dv2 [25] is a 3D hand-object dataset that contains 66,034 training samples and 11,524 evaluation samples. The annotations of the evaluation set are not available, so we use the official server for online evaluation. We also use this dataset to evaluate temporal performance.

We use the following metrics in quantitative evaluations.

MPJPE/MPVPE measures the mean per joint/vertex position error by Euclidean distance (mm) between the estimated and ground-truth coordinates.

PA-MPJPE/MPVPE is a modification of MPJPE/MPVPE with Procrustes analysis [21], ignoring global variation. For conciseness, this metric is abbreviated as PJ/PV.

Acc captures the acceleration of 2D/3D joint landmarks in pixel/s² or mm/s² to reflect temporal coherence.

AUC means the area under the curve of PCK (percentage of correct keypoints) *vs.* error thresholds of $n \sim 50$ mm (for 3D measurement) or 0~30 pixel (for 2D measurement). According to [93, 92], $n = 0$ or 20.

F-Score is the harmonic mean between recall and precision between two meshes *w.r.t.* a specific distance threshold.

F@5/F@15 corresponds to a threshold of 5mm/15mm.

Mult-Adds counts multiply-add operations.

#Param indicates the number of parameters.

4.3. Ablation studies

Our stacked networks and complement data. We use CMR [11] as the baseline, which adopts FC for 2D-to-3D feature mapping and SpiralConv++ for 3D decoding. With the same hyperparameters, we only change the 2D encoding network and data setting for comparisons. On the one hand, we study different stacked structures with from-scratch training. As shown in Table 1, we use ResNet [30]

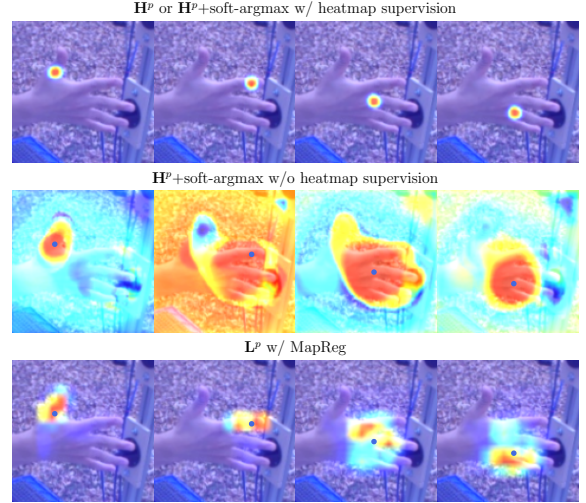


Figure 6. Visualization of *maps* in various 2D pose representations. Different from heatmap that focus on individual landmark, the *map* in MapReg adaptively describes inter-landmark constraints.

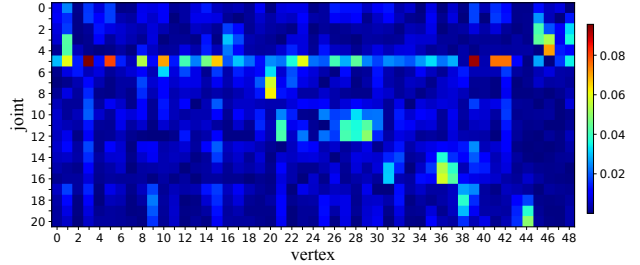


Figure 7. Learned lifting matrix.

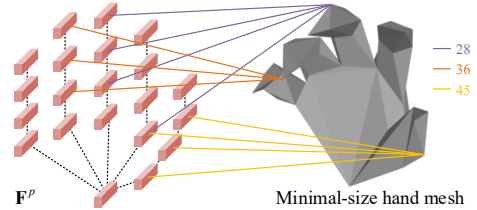


Figure 8. Highly relevant connections for several vertices in the lifting matrix. The numbers are vertex indices.

and MobileNet [67] to design stacked structures, the former of which contains intensive computation costs. Although MobileNet is computational tractability, it whittles down the performance (*i.e.*, PAMPJPE) by a large margin. On the other hand, ImageNet [15] is used for classification task, the knowledge from which is hard to transfer towards 2D/3D position regression. Consequently, the ImageNet pre-training brings less than 1mm PAMPJPE improvement.

Referring to Table 1, our complement data can be included in both pre-training and fine-tuning steps, and DenseStack/GhostStack induces 7.55/8.89mm PAMPJPE. Hence, we significantly reduce the computational cost without sacrificing reconstruction accuracy with our tai-

2D representation	2D part				2D-to-3D part							
	Mult-Adds	#Param	2D AUC↑	2D Acc↓	Pose pooling	Lifting	Mult-Adds	#Param	PJ↓	3D AUC↑	3D Acc↓	
—	—				Global pooling	FC	3.2M	3.2M	7.55	0.850	10.12	
H _p w/o skip.	57.7M	1.4M	0.796	2.57	Joint-wise pooling	PVL	0.4M	1029	7.41	0.852	9.11	
H _p	57.7M	1.4M	0.857	2.41		Grid sampling	PVL	0.3M	1029	7.28	0.856	8.32
H _p +soft-argmax	57.7M	1.4M	0.865	2.32		Grid sampling	PVL	0.3M	1029	7.11	0.859	6.14
L _p w/ reg.	0.1M	0.1M	0.847	2.16		Grid sampling	PVL	0.3M	1029	6.95	0.862	5.41
L _p w/ MapReg (ours)	47.5M	1.4M	0.870	2.04								

Table 2. Ablation study of feature lifting module. The first row is the setting of CMR [11] (*i.e.* the last row of Table 1); skip. and reg. denote skip connection and direct regression; 2D and 3D accuracy are tested on RHD and FreiHAND, respectively; Acc is *w.r.t.* HO3Dv2.

\mathcal{L}_{con3D}	\mathcal{L}_{con2D}	PJ \downarrow	3D AUC \uparrow	2D Acc \downarrow	3D Acc \downarrow
		6.95	0.862	2.04	5.41
✓		6.85	0.864	2.03	4.78
✓	✓	6.85	0.864	1.98	4.75

Table 3. Ablation study of consistency learning. The accuracy is *w.r.t.* FreiHAND and Acc is *w.r.t.* HO3Dv2.

lored encoding structures and dataset, making our DenseStack/GhostStack suitable for mobile environments.

Feature lifting module. With DenseStack, we perform a joint study of accuracy and temporal coherence. In Table 2, we do not use sequential module, temporal optimization, or post-processing. First, we explore various 2D pose representations in detail. \mathbf{H}^p (Figure 4(a)) is a high-accuracy representation, and the skip connection proves critical to fuse high- and low-resolution features. Hence, when the skip connection is removed, \mathbf{H}^p performs poorly in accuracy. As a differentiable form of picking maximum position, soft-argmax can produce smoother 2D positions from \mathbf{H}^p , hence improving both 2D AUC and Acc. In addition, \mathbf{L}^p w/ reg. (Figure 4(c)) has relatively modest accuracy but produces better temporal performance because of global receptive field. Ultimately, our MagReg achieves better 2D accuracy and temporal coherence by integrating the merits of heatmap- and regression-based paradigms.

To clearly reveal details, we illustrate *maps* in Figure 6. As for \mathbf{H}^p +soft-argmax, we additionally use the same loss setting (Equation 9) for training, *i.e.*, heatmap supervision is not involved. This manner naively induces a smooth version of \mathbf{H}^p because the heuristic soft-argmax neglects visual semantics. In terms of MapReg, we present the *map* before it is flattened into a vector. Different from \mathbf{H}^p , the *map* can adaptively describe joint landmark constraints, and then the 2D positions are predicted using adaptive local-global information. For example, the entire thumb is activated when predicting a landmark on it. Hence, our MapReg can produce more reasonable articulated structures (as shown in the supplemental material) to improve temporal coherence.

Subsequently, we explore 3D performance with 2D pose-aligned features. During pose pooling, joint-wise pooling [83] (Figure 4(e)) can be used to obtain 2D pose-aligned features based on \mathbf{H}^p , while grid sampling [66] (Figure 4(f)) is commonly adopted when soft-argmax or \mathbf{L}^p is utilized. It is noteworthy that although \mathbf{L}^p w/ reg. lags be-

3D decoding	Mult-Adds	#Param	PJ \downarrow	Acc \downarrow	FPS \uparrow
<i>w/ GhostStack</i>					
SpiralConv++	159.0/263.1M	1.0/6.2M	8.63	2.31/7.07	77
DSConv (ours)	19.5/123.5M	0.1/5.3M	8.76	2.30/6.98	83
<i>w/ DenseStack</i>					
SpiralConv++	159.0/579.4M	1.0/9.0M	6.85	1.98/4.75	59
DSConv (ours)	19.5/439.9M	0.1/8.1M	6.87	1.92/4.73	67

Table 4. Ablation study of 3D decoding. Mult-Adds and #Param are *w.r.t.* the 3D decoder/overall model; 2D/3D Acc is presented; FPS is tested on Apple A14 CPU; Accuracy and temporal performance are tested on FreiHAND and HO3Dv2, respectively.

hind in terms of 2D accuracy, it produces better PAMPJPE than \mathbf{H}^p +soft-argmax. Hence, with the same pose pooling method, 2D consistency is more crucial than accuracy in establishing a stable training process. Ultimately, MagRep-based \mathbf{L}^p induces the best PAMPJPE and 3D Acc.

In PVL, we design a linear operation with lifting matrix $\mathbf{M}^l \in \mathbb{R}^{V^{mini} \times N}$ to transform features from 2D pose space to 3D vertex space. Thus, V^{mini} vertex features are produced by a linear combination of N landmark features. Figure 7 depicts a well-trained lifting matrix, where we illustrate $\text{abs}(\mathbf{M}^l)$ to clearly reveal the joint-vertex relations. As can be observed, the learned \mathbf{M}^l is sparse. A joint landmark (*i.e.*, joint 5 that locates at the root of forefinger) serves as the global information and contributes to the majority of vertices. Besides, some joint landmark traits are propagated to their corresponding vertices. Figure 8 depicts highly relevant connections in \mathbf{M}^l , demonstrating that the PVL approach can preserve the semantic consistency.

When compared to CMR [11] (the first row of Table 2), our feature lifting module results in better PAMPJPE and 3D Acc. Moreover, we significantly reduce the computational expenses in the 2D-to-3D part. Despite the use of extra Mult-Adds in the 2D part, 2D pose prediction has previously proven to be beneficial in 3D hand reconstruction due to multi-task learning [85] and root recovery task [11].

Towards balancing model efficiency and performance.

We design 3D/2D consistency loss to improve the performance further. From Table 3, we can see that consistency learning improves temporal coherence. Besides, the accuracy and temporal coherence can benefit from each other so that PAMPJPE is also enhanced.

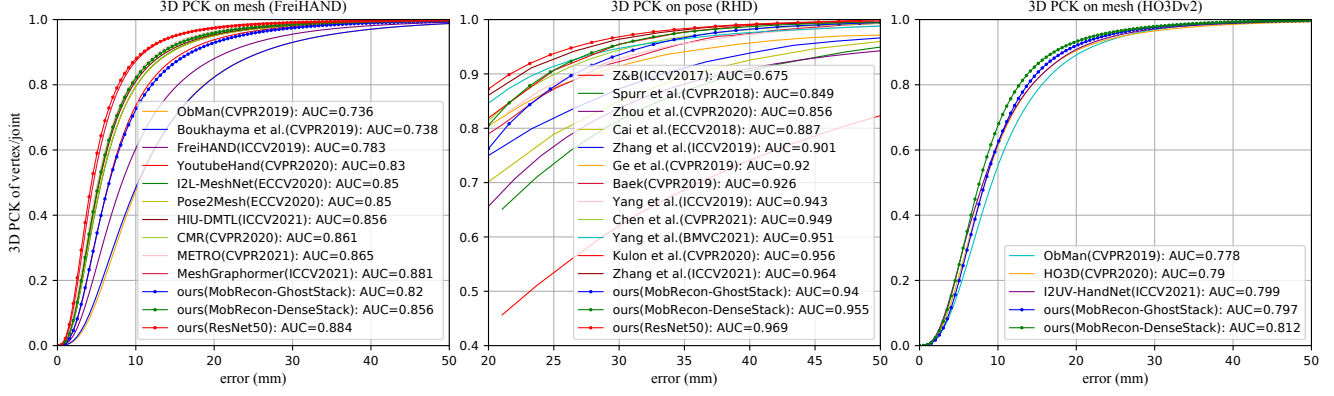


Figure 9. 3D PCK vs. error thresholds.

Method	Backbone	PJ ↓	PV ↓	F@5 ↑	F@15 ↑
MobileHand [22]	MobileNet	—	13.1	0.439	0.902
FreiHAND [93]	ResNet50	11.0	10.9	0.516	0.934
YotubeHand [46]	ResNet50	8.4	8.6	0.614	0.966
I2L-MeshNet [57]	ResNet50*	7.4	7.6	0.681	0.973
HIU-DMTL [85]	Customized*	7.1	7.3	0.699	0.974
CMR [11]	ResNet50*	6.9	7.0	0.715	0.977
I2UV-HandNet [10]	ResNet50	6.7	6.9	0.707	0.977
METRO [50]	HRNet	6.7	6.8	0.717	0.981
Tang <i>et al.</i> [73]	ResNet50	6.7	6.7	0.724	0.981
MeshGraphormer [51]	HRNet	5.9	6.0	0.765	0.987
MobRecon (ours)	GhostStack*	8.8	9.1	0.597	0.960
MobRecon (ours)	DenseStack*	6.9	7.2	0.694	0.979
ours [‡]	ResNet18*	6.7	6.8	0.727	0.979
ours [†]	ResNet18*	6.1	6.3	0.758	0.983
ours [‡]	ResNet50*	6.1	6.2	0.760	0.984
ours [†]	ResNet50*	5.7	5.8	0.784	0.986

Table 5. Results on the FreiHAND dataset. *: stacked structure;

†: These models are based on ImageNet pre-trained backbone and mixed fine-tuning data; ‡: These models are totally unrelated to our complement data.

As shown in Table 4, DSConv dramatically decreases the Mult-Adds and #Param of the 3D decoder and obtains on par, sometimes even better, performance compared with SpiralConv++. Overall, our MobRecon with DenseStack/GhostStack can reach 67/83 FPS on Apple A14 CPU.

Discussion. MobRecon has a limitation that the DSConv increases memory access cost, so some engineering optimization should be involved for higher inference speed.

4.4. Comparisons with Contemporary Methods

On the FreiHAND dataset, we scale up our ResNet-based model with 224×224 input resolution for a fair comparison. As shown in Table 5, we surpass previous methods with ResNet50, leading to a new state of the art, *i.e.*, 5.7mm PAMPJPE. Based on DenseStack or GhostStack, our MobRecon outmatches some ResNet-based methods. Referring to Figure 9, the proposed MobRecon has superior performance on 3D PCK. Beyond high accuracy, we also

Method	Backbone	PJ ↓	PV ↓	F@5 ↑	F@15 ↑
ObMan [29]	ResNet18	11.0	11.0	0.464	0.939
HO3D [25]	CPM [77]	10.7	10.6	0.506	0.942
I2UV-HandNet [10]	ResNet50	9.9	10.1	0.500	0.943
MobRecon (ours)	GhostStack	10.0	10.2	0.488	0.948
MobRecon (ours)	DenseStack	9.2	9.4	0.538	0.957

Table 6. Results on the HO3Dv2 dataset.

achieve superior inference speed, as shown in Figure 1.

In experiments on RHD and HO3Dv2, our complement data are only used to pre-train DenseStack/GhostStack. On the RHD dataset, we compare with several pose estimation methods (such as [1, 5, 69, 81]) in Figure 9. Our MobRecon with DenseStack/GhostStack induces 3D AUC of 0.955 and 0.940, outperforming most compared approaches.

The HO3Dv2 dataset is employed for evaluation. According to Table 6, our MobRecon outperforms existing methods such as [29, 25, 10]. HO3Dv2 is more challenging than FreiHAND and RHD because of serious object occlusion. In this case, MobRecon outperforms some ResNet-based methods because of better generalization ability. Besides, we also achieve better temporal coherence in this sequential task, as shown in Table 2.

Please refer to the supplementary materials for more qualitative analyses, mobile applications, *etc.*

5. Conclusions and Future Work

In this work, we present a novel hand mesh reconstruction method with superior efficiency, accuracy, and temporal coherence. First, we propose lightweight stacked structures for 2D encoding. Then, a feature lifting module with MapReg, pose pooling, and PVL approaches is designed for 2D-to-3D mapping. Besides, DSConv is developed to handle the 3D decoding task efficiently. Our MobRecon only involves 123M Mult-Adds and 5M parameters so that it reaches a fast inference speed of 83FPS on Apple A14 CPU. Moreover, we achieve the state-of-the-art performance on FreiHAND, RHD, and HO3Dv2. We plan to investigate efficient methods for interacting hands.

References

- [1] Seungryul Baek, Kwang In Kim, and Tae-Kyun Kim. Pushing the envelope for RGB-based dense 3D hand pose estimation via neural rendering. In *CVPR*, 2019. 2, 8
- [2] Seungryul Baek, Kwang In Kim, and Tae-Kyun Kim. Weakly-supervised domain adaptation via GAN and mesh model for estimating 3D hand poses interacting objects. In *CVPR*, 2020. 2
- [3] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. LoopReg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. In *NeurIPS*, 2020. 2
- [4] Adnane Boukhayma, Rodrigo de Bem, and Philip HS Torr. 3D hand shape and pose from images in the wild. In *CVPR*, 2019. 2, 12
- [5] Yujun Cai, Lihao Ge, Jianfei Cai, and Junsong Yuan. Weakly-supervised 3D hand pose estimation from monocular RGB images. In *ECCV*, 2018. 8
- [6] Zhe Cao, Ilija Radosavovic, Angjoo Kanazawa, and Jitendra Malik. Reconstructing hand-object interactions in the wild. In *ICCV*, 2021. 2
- [7] Joao Carreira, Pulkit Agrawal, Katerina Fragkiadaki, and Jitendra Malik. Human pose estimation with iterative error feedback. In *CVPR*, 2016. 4
- [8] Yu-Wei Chao, Wei Yang, Yu Xiang, Pavlo Molchanov, Ankur Handa, Jonathan Tremblay, Yashraj S. Narang, Karl Van Wyk, Umar Iqbal, Stan Birchfield, Jan Kautz, and Dieter Fox. DexYCB: A benchmark for capturing hand grasping of objects. In *CVPR*, 2021. 12
- [9] Liangjian Chen, Shih-Yao Lin, Yusheng Xie, Yen-Yu Lin, and Xiaohui Xie. MVHM: A large-scale multi-view hand mesh benchmark for accurate 3D hand pose estimation. In *WACV*, 2021. 12
- [10] Ping Chen, Yujin Chen, Dong Yang, Fangyin Wu, Qin Li, Qingpei Xia, and Yong Tan. I2uv-handnet: Image-to-uv prediction network for accurate and high-fidelity 3d hand mesh modeling. In *ICCV*, pages 12929–12938, 2021. 2, 8
- [11] Xingyu Chen, Yufeng Liu, Chongyang Ma, Jianlong Chang, Huayan Wang, Tian Chen, Xiaoyan Guo, Pengfei Wan, and Wen Zheng. Camera-space hand mesh recovery via semantic aggregation and adaptive 2D-1D registration. In *CVPR*, 2021. 1, 2, 3, 4, 6, 7, 8, 14
- [12] Yujin Chen, Zhigang Tu, Di Kang, Linchao Bao, Ying Zhang, Xuefei Zhe, Ruizhi Chen, and Junsong Yuan. Model-based 3d hand reconstruction via self-supervised learning. In *CVPR*, 2021. 2
- [13] Hongsuk Choi, Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. Beyond static features for temporally consistent 3D human pose and shape from a video. In *CVPR*, 2021. 1, 2, 3
- [14] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NeurIPS Workshop*, 2014. 2
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009. 6
- [16] Guillermo Garcia-Hernando, Shanxin Yuan, Seungryul Baek, and Tae-Kyun Kim. First-person hand action benchmark with rgb-d videos and 3d hand pose annotations. In *CVPR*, 2018. 12
- [17] Michael Garland and Paul S Heckbert. Surface simplification using quadric error metrics. In *Proceedings of the Annual Conference on Computer Graphics and Interactive Techniques*, 1997. 4
- [18] Lihao Ge, Zhou Ren, Yuncheng Li, Zehao Xue, Yingying Wang, Jianfei Cai, and Junsong Yuan. 3D hand shape and pose estimation from a single RGB image. In *CVPR*, 2019. 1, 2, 3, 4, 5
- [19] Francisco Gomez-Donoso, Sergio Orts-Escolano, and Miguel Cazorla. Large-scale multiview 3D hand pose dataset. *Image and Vision Computing*, 2019. 12
- [20] Shunwang Gong, Lei Chen, Michael Bronstein, and Stefanos Zafeiriou. SpiralNet++: A fast and highly efficient mesh convolution operator. In *ICCV Workshops*, 2019. 5
- [21] John C Gower. Generalized procrustes analysis. *Psychometrika*, 1975. 6
- [22] Lim Guan Ming, Jatesiktat Prayook, and Ang Wei Tech. MobileHand: Real-time 3D hand shape and pose estimation from color image. In *ICONIP*, 2020. 1, 2, 8
- [23] Riza Alp Guler and Iasonas Kokkinos. HoloPose: Holistic 3D human reconstruction in-the-wild. In *CVPR*, 2019. 3, 4
- [24] Jia Guo, Jiankang Deng, Niannan Xue, and Stefanos Zafeiriou. Stacked dense U-Nets with dual transformers for robust face alignment. In *BMVC*, 2018. 5
- [25] Shreyas Hampali, Mahdi Rad, Markus Oberweger, and Vincent Lepetit. Honnotate: A method for 3d annotation of hand and object poses. In *CVPR*, 2020. 5, 6, 8, 12
- [26] Shreyas Hampali, Sayan Deb Sarkar, Mahdi Rad, and Vincent Lepetit. HandsFormer: Keypoint transformer for monocular 3d pose estimation of hands and object in interaction. *arXiv:2104.14639*, 2021. 3
- [27] Kai Han, Yunhe Wang, Qi Tian, Jianyuan Guo, Chunjing Xu, and Chang Xu. GhostNet: More features from cheap operations. In *CVPR*, 2020. 2, 3
- [28] Shangchen Han, Beibei Liu, Randi Cabezas, Christopher D Twigg, Peizhao Zhang, Jeff Petkau, Tsz-Ho Yu, Chun-Jung Tai, Muzaffer Akbay, Zheng Wang, et al. MEgATrack: monochrome egocentric articulated hand-tracking for virtual reality. *TOG*, 2020. 1
- [29] Yana Hasson, Gul Varol, Dimitrios Tzionas, Igor Kaleyvtykh, Michael J Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *CVPR*, 2019. 2, 8, 12
- [30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 2, 6
- [31] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997. 4
- [32] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu,

- Ruoming Pang, Vijay Vasudevan, et al. Searching for MobileNetv3. In *ICCV*, 2019. 1, 2
- [33] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv:1704.04861*, 2017. 2, 5
- [34] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018. 3
- [35] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, 2017. 3
- [36] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. ARCH: Animatable reconstruction of clothed humans. In *CVPR*, 2020. 2
- [37] Umar Iqbal, Pavlo Molchanov, Thomas Breuel Juergen Gall, and Jan Kautz. Hand pose estimation via latent 2.5D heatmap regression. In *ECCV*, 2018. 2
- [38] Timothy Jeruzalski, Boyang Deng, Mohammad Norouzi, John P Lewis, Geoffrey Hinton, and Andrea Tagliasacchi. NASA: Neural articulated shape approximation. In *ECCV*, 2020. 2
- [39] Hanwen Jiang, Shaowei Liu, Jiashun Wang, and Xiaolong Wang. Hand-object contact consistency reasoning for human grasps generation. In *ICCV*, 2021. 2
- [40] Hanbyul Joo, Tomas Simon, Mina Cikara, and Yaser Sheikh. Towards social artificial intelligence: Nonverbal social signal prediction in a triadic interaction. In *CVPR*, 2019. 1
- [41] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3D human dynamics from video. In *CVPR*, 2019. 1, 2
- [42] Korrawe Karunratanakul, Adrian Spurr, Zicong Fan, Otmar Hilliges, and Siyu Tang. A skeleton-driven neural occupancy representation for articulated hands. *arXiv:2109.11399*, 2021. 2
- [43] Korrawe Karunratanakul, Jinlong Yang, Yan Zhang, Michael J Black, Krikamol Muandet, and Siyu Tang. Grasping field: Learning implicit representations for human grasps. In *3DV*, 2020. 2
- [44] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014. 5
- [45] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. VIBE: Video inference for human body pose and shape estimation. In *CVPR*, 2020. 1, 2
- [46] Dominik Kulon, Riza Alp Guler, Iasonas Kokkinos, Michael M Bronstein, and Stefanos Zafeiriou. Weakly-supervised mesh-convolutional hand reconstruction in the wild. In *CVPR*, 2020. 1, 2, 3, 4, 5, 8, 12
- [47] Taein Kwon, Bugra Tekin, Jan Stuhmer, Federica Bogo, and Marc Pollefeys. H2o: Two hands manipulating objects for first person interaction recognition. In *ICCV*, 2021. 12
- [48] Jiefeng Li, Siyuan Bian, Ailing Zeng, Can Wang, Bo Pang, Wentao Liu, and Cewu Lu. Human pose regression with residual log-likelihood estimation. In *ICCV*, 2021. 3, 4
- [49] Isaak Lim, Alexander Dielen, Marcel Campen, and Leif Kobbelt. A simple approach to intrinsic correspondence learning on unstructured 3D meshes. In *ECCV*, 2018. 2, 4
- [50] Kevin Lin, Lijuan Wang, and Zicheng Liu. End-to-end human pose and mesh reconstruction with transformers. In *CVPR*, 2021. 1, 2, 8
- [51] Kevin Lin, Lijuan Wang, and Zicheng Liu. Mesh graphormer. In *ICCV*, 2021. 1, 2, 8
- [52] Shaowei Liu, Hanwen Jiang, Jiarui Xu, Sifei Liu, and Xiaolong Wang. Semi-supervised 3D hand-object poses estimation with interactions in time. In *CVPR*, 2021. 2
- [53] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *TOG*, 2015. 3
- [54] Zhengyi Luo, S Alireza Golestaneh, and Kris M Kitani. 3D human motion estimation via motion compression and refinement. In *ACCV*, 2020. 1, 2
- [55] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *CVPR*, 2019. 2
- [56] Marko Mihajlovic, Yan Zhang, Michael J Black, and Siyu Tang. LEAP: Learning articulated occupancy of people. In *CVPR*, 2021. 2
- [57] Gyeongsik Moon and Kyoung Mu Lee. I2L-MeshNet: Image-to-voxel prediction network for accurate 3D human pose and mesh estimation from a single RGB image. In *ECCV*, 2020. 2, 8
- [58] Gyeongsik Moon and Kyoung Mu Lee. Pose2Pose: 3d positional pose-guided 3d rotational pose prediction for expressive 3D human pose and mesh estimation. *arXiv:2011.11534*, 2020. 3
- [59] Gyeongsik Moon, Shoou-I Yu, He Wen, Takaaki Shiratori, and Kyoung Mu Lee. Interhand2.6M: A dataset and baseline for 3d interacting hand pose estimation from a single RGB image. In *ECCV*, 2020. 2, 12
- [60] Franziska Mueller, Florian Bernard, Oleksandr Sotnychenko, Dushyant Mehta, Srinath Sridhar, Dan Casas, and Christian Theobalt. Generated hands for real-time 3d hand tracking from monocular rgb. In *CVPR*, 2018. 12
- [61] Franziska Mueller, Dushyant Mehta, Oleksandr Sotnychenko, Srinath Sridhar, Dan Casas, and Christian Theobalt. Real-time hand tracking under occlusion from an egocentric RGB-D sensor. In *ICCV*, 2017. 12
- [62] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016. 2, 3, 4
- [63] Evonne Ng, Shiry Ginosar, Trevor Darrell, and Hanbyul Joo. Body2Hands: Learning to infer 3D hands from conversational gesture body dynamics. In *CVPR*, 2021. 1
- [64] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *CVPR*, 2021. 2
- [65] Javier Romero, Dimitrios Tzionas, and Michael J Black. Embodied hands: Modeling and capturing hands and bodies together. *TOG*, 2017. 2, 4
- [66] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. PIFu: Pixel-

- aligned implicit function for high-resolution clothed human digitization. In *ICCV*, 2019. 2, 3, 4, 7
- [67] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *CVPR*, 2018. 2, 6
- [68] Adrian Spurr, Aneesh Dahiya, Xi Wang, Xucong Zhang, and Otmar Hilliges. Self-supervised 3d hand pose estimation from monocular rgb via contrastive learning. In *ICCV*, 2021. 5
- [69] Adrian Spurr, Jie Song, Seonwook Park, and Otmar Hilliges. Cross-modal deep variational hand pose estimation. In *CVPR*, 2018. 8
- [70] Srinath Sridhar, Franziska Mueller, Michael Zollhöfer, Dan Casas, Antti Oulasvirta, and Christian Theobalt. Real-time joint tracking of a hand manipulating an object from RGB-D input. In *ECCV*, 2016. 12
- [71] Xiao Sun, Jiaxiang Shang, Shuang Liang, and Yichen Wei. Compositional human pose regression. In *ICCV*, 2017. 4
- [72] Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. Integral human pose regression. In *ECCV*, 2018. 4
- [73] Xiao Tang, Tianyu Wang, and Chi-Wing Fu. Towards accurate alignment in real-time 3d hand-mesh reconstruction. In *ICCV*, 2021. 1, 8
- [74] Alexander Toshev and Christian Szegedy. DeepPose: Human pose estimation via deep neural networks. In *CVPR*, 2014. 4
- [75] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. Deep high-resolution representation learning for visual recognition. *TPAMI*, 2020. 1, 4
- [76] Fangyun Wei, Xiao Sun, Hongyang Li, Jingdong Wang, and Stephen Lin. Point-set anchors for object detection, instance segmentation and pose estimation. In *ECCV*, 2020. 4
- [77] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *CVPR*, 2016. 8
- [78] John Yang, Hyung Jin Chang, Seungeui Lee, and Nojun Kwak. SeqHAND: RGB-sequence-based 3D hand pose and shape estimation. In *European Conference on Computer Vision*, pages 122–139. Springer, 2020. 12
- [79] Linlin Yang, Shicheng Chen, and Angela Yao. SemiHand: Semi-supervised hand pose estimation with consistency. In *ICCV*, 2021. 2, 5
- [80] Lixin Yang, Jiasen Li, Wenqiang Xu, Yiqun Diao, and Cewu Lu. BiHand: Recovering hand mesh with multi-stage bisected hourglass networks. In *BMVC*, 2020. 2, 3
- [81] Linlin Yang, Shile Li, Dongheui Lee, and Angela Yao. Aligning latent spaces for 3D hand pose estimation. In *ICCV*, 2019. 8
- [82] Lixin Yang, Xinyu Zhan, Kailin Li, Wenqiang Xu, Jiefeng Li, and Cewu Lu. CPF: Learning a contact potential field to model the hand-object interaction. In *ICCV*, 2021. 2
- [83] Baowen Zhang, Yangang Wang, Xiaoming Deng, Yinda Zhang, Ping Tan, Cuixia Ma, and Hongan Wang. Interacting two-hand 3D pose and shape reconstruction from single color image. In *ICCV*, 2021. 2, 3, 4, 7
- [84] Jiawei Zhang, Jianbo Jiao, Mingliang Chen, Liangqiong Qu, Xiaobin Xu, and Qingxiong Yang. 3D hand pose tracking and estimation using stereo matching. *arXiv:1610.07214*, 2016. 12
- [85] Xiong Zhang, Hongsheng Huang, Jianchao Tan, Hongmin Xu, Cheng Yang, Guozhu Peng, Lei Wang, and Ji Liu. Hand image understanding via deep multi-task learning. In *ICCV*, 2021. 1, 2, 4, 7, 8
- [86] Xiong Zhang, Qiang Li, Hong Mo, Wenbo Zhang, and Wen Zheng. End-to-end hand mesh recovery from a monocular RGB image. In *ICCV*, 2019. 2
- [87] Zhengyi Zhao, Tianyao Wang, Siyu Xia, and Yangang Wang. Hand-3D-Studio: A new multi-view system for 3D hand reconstruction. In *ICASSP*, 2020. 12
- [88] Zimeng Zhao, Xi Zhao, and Yangang Wang. TravelNet: Self-supervised physically plausible hand motion learning from monocular color images. In *ICCV*, 2021. 2
- [89] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. PaMIR: Parametric model-conditioned implicit representation for image-based human reconstruction. *TPAMI*, 2021. 3
- [90] Yuxiao Zhou, Marc Habermann, Weipeng Xu, Ikhsanul Habibie, Christian Theobalt, and Feng Xu. Monocular real-time hand shape and motion capture using multi-modal data. In *CVPR*, 2020. 1, 2
- [91] Christian Zimmermann, Max Argus, and Thomas Brox. Contrastive representation learning for hand shape estimation. *arXiv:2106.04324*, 2021. 5
- [92] Christian Zimmermann and Thomas Brox. Learning to estimate 3D hand pose from single RGB images. In *ICCV*, 2017. 6, 12
- [93] Christian Zimmermann, Duygu Ceylan, Jimei Yang, Bryan Russell, Max Argus, and Thomas Brox. FreiHAND: A dataset for markerless capture of hand pose and shape from single RGB images. In *ICCV*, 2019. 2, 5, 6, 8, 12

Dataset	Type	Size	Mesh	UP	MV
STB [84]	real	36K	×	×	×
RHD [92]	synthetic	44K	×	×	×
GANerated Hands [60]	synthetic	331K	×	×	×
SeqHAND [78]	synthetic	410K	✓	×	×
EgoDexter [61]	real	3K	×	×	×
Dexter+Object [70]	real	3K	×	×	×
FreiHAND [93]	real	134K	✓	×	×
YoutubeHand [46]	real	47K	✓	×	×
ObMan [29]	synthetic	153K	✓	×	×
HO3D [25]	real	77K	✓	×	×
DexYCB [8]	real	528K	✓	×	×
H2O [47]	synthetic	571K	✓	×	×
FPFA [16]	synthetic	105K	×	×	×
H3D [87]	real	22K	✓	×	✓(15)
MHP [19]	real	80K	×	×	✓(4)
MVHM [9]	synthetic	320K	✓	×	✓(8)
InterHand2.6M [59]	real	2.6M	✓	×	✓(80)
ours	synthetic	328K	✓	✓	✓(216)

Table 7. Comparison among RGB-based 3D hand datasets. “MV” means multi-view, and the number in brackets shows the total number of views. “UP” denotes uniform pose distribution.

6. Complement Dataset

Motivation. As shown in Table 7, many datasets are developed for 3D hand pose estimation [84, 92, 61, 70, 93, 46, 19, 87, 8, 47, 16, 78, 60]. To collect real-world hand data, existing datasets are usually captured using a multi-view studio and annotated via semi-automatic model fitting [93, 19]. However, these model-fitted datasets usually suffer from noisy annotation, lack of background diversity, and costly data collection. An alternative way is computer-aided synthetic data [92, 9], which are superior in scalability, distribution, annotation, and collection cost. In addition, a good training dataset should avoid long-tailed distributions. That is, both hand poses and viewpoints should be uniformly distributed. Unfortunately, we are not aware of any existing dataset that fits this requirement. Some datasets try to alleviate the problem of limited viewpoints by multi-view rendering (e.g., MVHM [9] contains 8 views), but they are still too sparse to cover all the possible views. Boukhayma *et al.* [4] uniformly sampled MANO PCA components to produce various hand poses. However, the PCA space does not describe physical factors, so the corresponding sampling results cannot be intuitively controlled. Thereby, we are inspired to generate a more comprehensive hand dataset with sufficient and uniformly distributed hand poses and viewpoints.

Data Designs. We design a high-fidelity hand mesh with 5633 vertices and 11232 faces. Different from existing hand datasets, we uniformly design hand poses. First, as shown in Figure 10, we set two states for each finger, *i.e.*, total bending and extending. Then, we obtain 32 base poses by combining five finger states. The combination of these



Figure 10. Base poses. Under the consideration of politeness, one pose with middle finger extending is not shown.

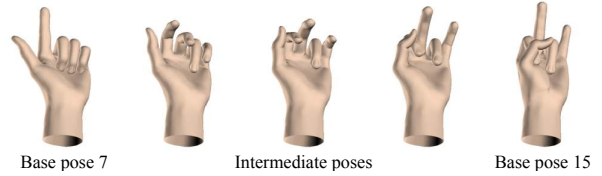


Figure 11. Intermediate poses from base pose 7 to base pose 15.

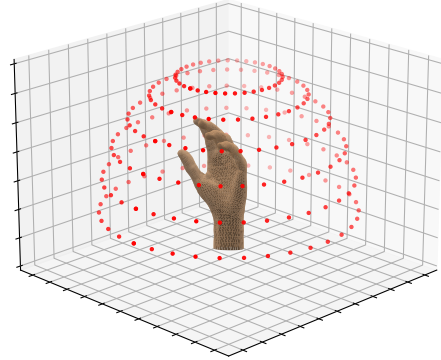


Figure 12. Illustration of viewpoints to render the dataset. Each red point denotes a camera position. The camera points to the palm center.

base poses results in 496 pose pairs. For each pair, we uniformly interpolate three intermediate poses from one pose to another in Maya software¹ (as shown in Figure 11). In total, we obtain 1520 uniformly distributed hand poses.

For each pose sample, we provide its dense viewpoints by rendering. To this end, we uniformly define 216 hemispherical-arranged camera positions. As shown in Figure 12, the longitude ranges from 0 to 2π while the latitude ranges from 0 to $\pi/2$. Adjacent positions differ in longitude or latitude by $\pi/18$ or $\pi/12$. All cameras point to the palm center so that the hand locates at the center of rendered images. Because the end of the wrist locates at the sphere center, the hemispherical sampling contains the

¹<https://www.autodesk.com/>

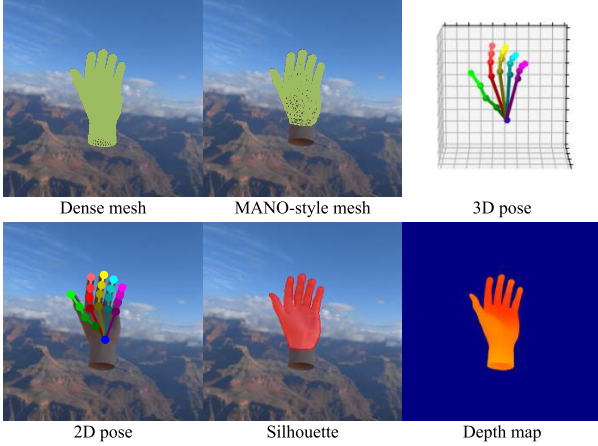


Figure 13. Annotations.

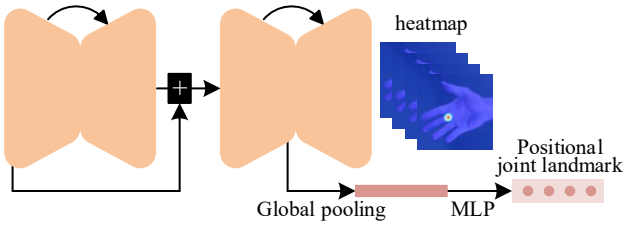


Figure 14. The pre-training architecture using our data.

first-person perspectives. As for the background, we collect high-dynamic-range (HDR) imaging with real scenes and illumination for rendering so that our hand mesh can realistically blend into various scenes. Figure 20 illustrates rendered samples with our viewpoints.

The automatically generated annotations involve no noise. Consistent with mainstream datasets, we design a pose-agnostic matrix to map our dense topology to MANO-style mesh with 778 vertices and 1538 faces. As shown in Figure 13, we provide annotations of our designed dense mesh, MANO-style mesh, 3D pose, 2D pose, silhouette, depth map, and intrinsic camera parameters.

Discussion. The limitation of our data is the lack of shape/texture diversity. Additionally, we only consider finger bend, and we plan to model finger splay to extend this dataset to cover the entire pose space uniformly.

Network pre-training. To pre-train the 2D encoding network, we design a 2D pose estimation task without the need of 3D annotation. In the main text, we analyze 2D representations with heatmap and position regression. Hence, as shown in Figure 14, we equally consider these representations during the pre-training step. That is, both heatmap and positional joint landmark are supervised. The model is pre-trained for 80 epochs with a mini-batch size of 128. The initial learning rate is 10^{-3} , which is divided by

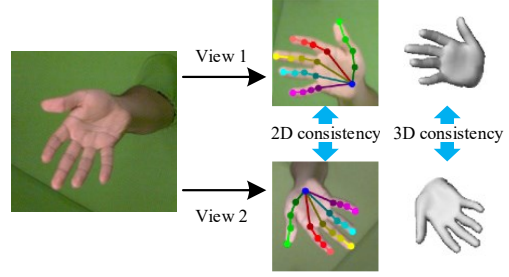


Figure 15. Consistency loss based on data augmentation.

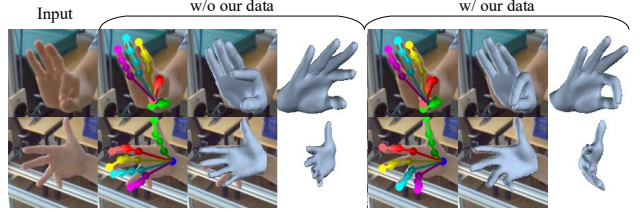


Figure 16. Qualitative visualization of 2D pose, aligned mesh, and side-view mesh.

10 at the 20th, 40th, and 60th epochs. The input resolution is 128×128 .

7. Analysis and Application

Diagram of our consistency loss. As shown in Figure 15, two views are derived with data augmentation with an input image. Then consistency loss can be designed in both 2D and 3D spaces as Equation 11 in the main text.

Explanation of dataset setting. During the ablation study in the main text, we use RHD, FreiHAND, and HO3Dv2 to evaluate different properties. Because FreiHAND and HO3Dv2 do not release ground truth and the official tools do not support 2D evaluation, RHD is employed for testing 2D accuracy. HO3Dv2 is a sequential dataset, so it is adopted to reflect temporal coherence. However, HO3Dv2 highlights hand-object interaction, which is not our topic. In contrast, FreiHAND highlights various hand poses, lighting conditions, *etc.*, so we use it for evaluating 3D accuracy.

The effect of our complement data during fine-tuning. As shown in Figure 16, when our data are employed in fine-tuning step, it can improve the model performance on difficult pose prediction.

Visualization on HO3Dv2. Referring to Table 6 in the main text, our MobRecon outperforms some ResNet-based models. We observe that this phenomenon is related to generalization performance. As shown in Figure 17, HO3Dv2 contains massive seriously occluded samples. Under this extreme condition, our model can produce a physically correct prediction while the ResNet-based model collapses.

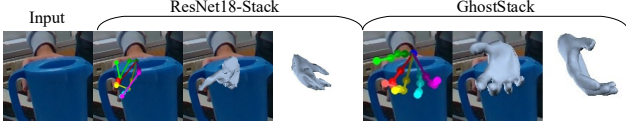


Figure 17. Comparison of GhostStack and ResNet18 on a challenging HO3Dv2 sample.

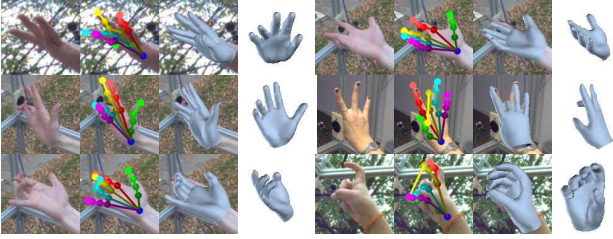


Figure 18. Typical failure cases.

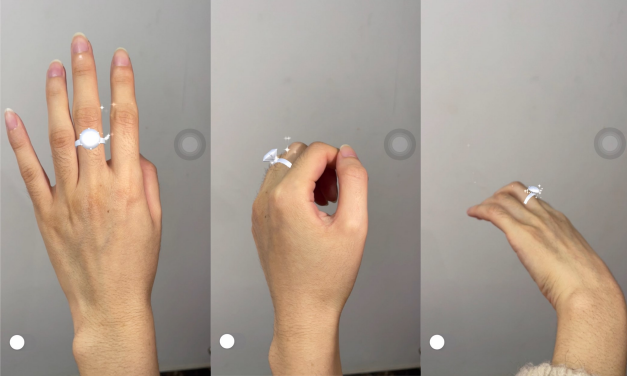


Figure 19. We develop an AR effect with MobRecon and deploy it on mobile devices. This figure is captured with iPhone12.

Failure case analysis. As shown in Figure 18, MobRecon could suffer from failure cases as for challenging poses. Typically, self-occlusion by finger splay is hard to accurately predict because they are tail-distributed poses in most datasets. We will solve this problem by improving our complement data, as stated in the above section.

More qualitative results. Figure 21 illustrates comprehensive qualitative results of our predicted 2D pose, aligned and side-view mesh. The challenges include challenging poses, object occlusion, truncation, and bad illumination. Overcoming these difficulties, our method can generate accurate 2D pose and 3D mesh.

Qualitative comparison on temporal coherence. We record a video snippet to demonstrate temporal coherence, where we keep the camera and hand static to produce low acceleration. Despite the static condition, the network input could be temporally unstable because of detection jitter *etc.* The ground-truth pose is straightforward (see Figure 22), and all compared models can easily obtain high accuracy. Hence, temporal performance can be exclusively

revealed in this experiment. As shown in Figure 22, our MobRecon performs better than CMR [11] in terms of 2D/3D pose consistency. In addition, we also compute the root coordinates with the method in [11] and achieve better root recovery stability. Besides, we also complement 2D PCK curves on RHD, which demonstrate that our method has better 2D pose accuracy. Beyond accuracy and temporal coherence, our MobRecon with MapReg can produce better articulated structures because of global receptive field and adaptive inter-landmark constraints (see Figure 6 in the main text).

Mobile application. Based on our MobRecon, a virtual ring can be worn with AR technique (Figure 19).

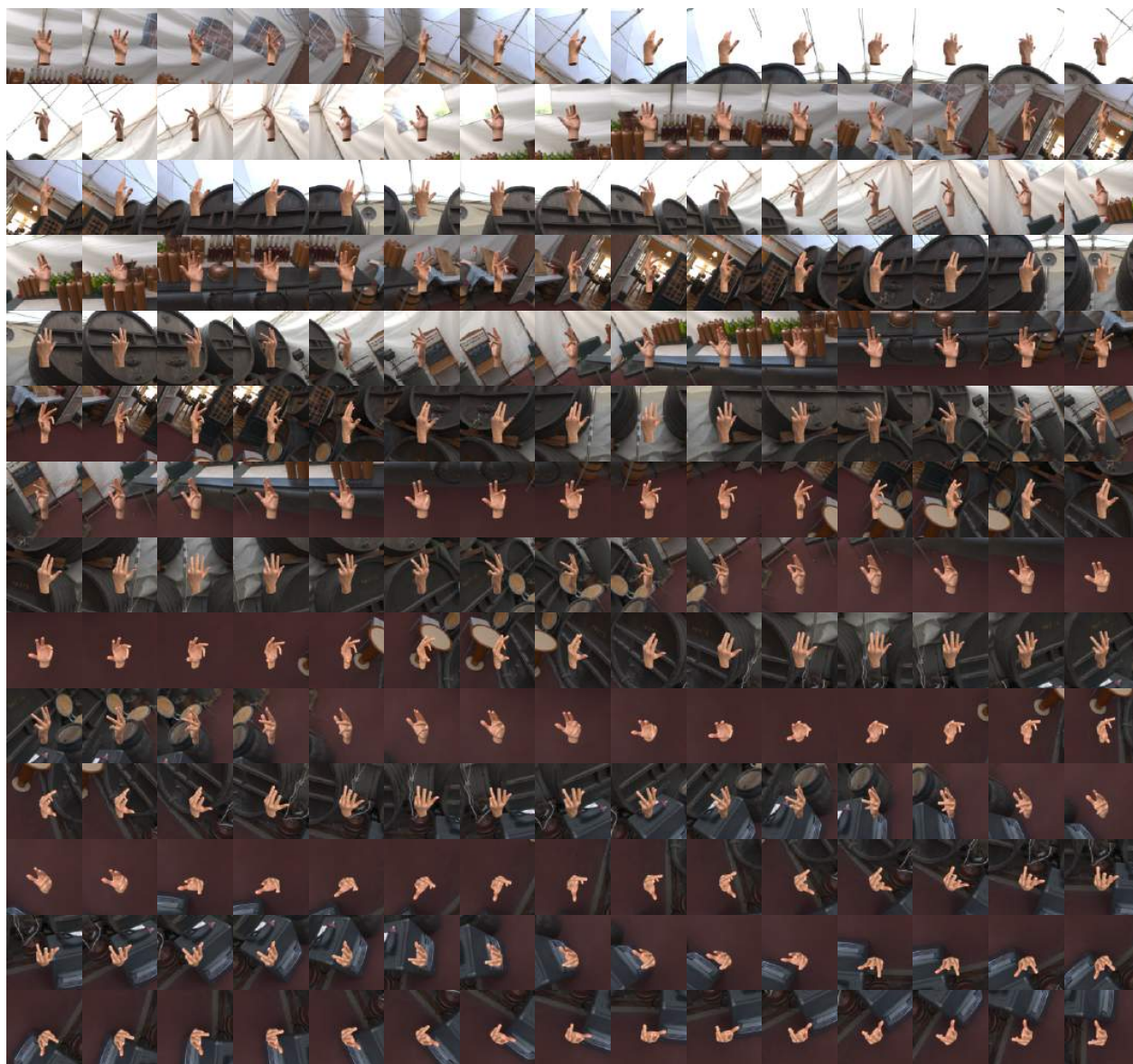


Figure 20. Rendered samples with dense and uniform viewpoints.

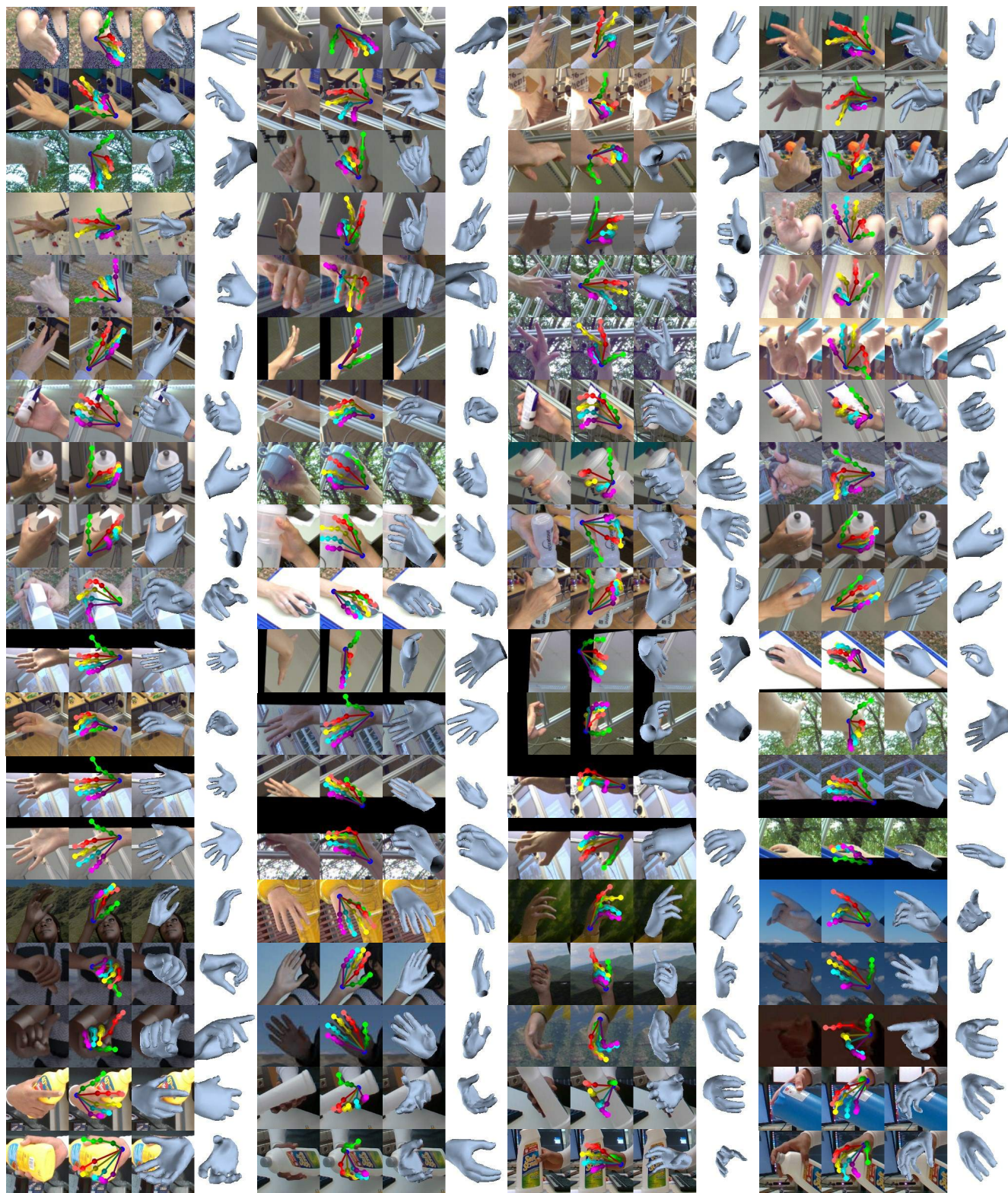


Figure 21. Qualitative visualization of 2D pose, aligned mesh, and side-view mesh on FreiHAND, RHD, and HO3Dv2. Our method is robust enough to handle cases of occlusion, truncation, challenging poses, and bad illumination.

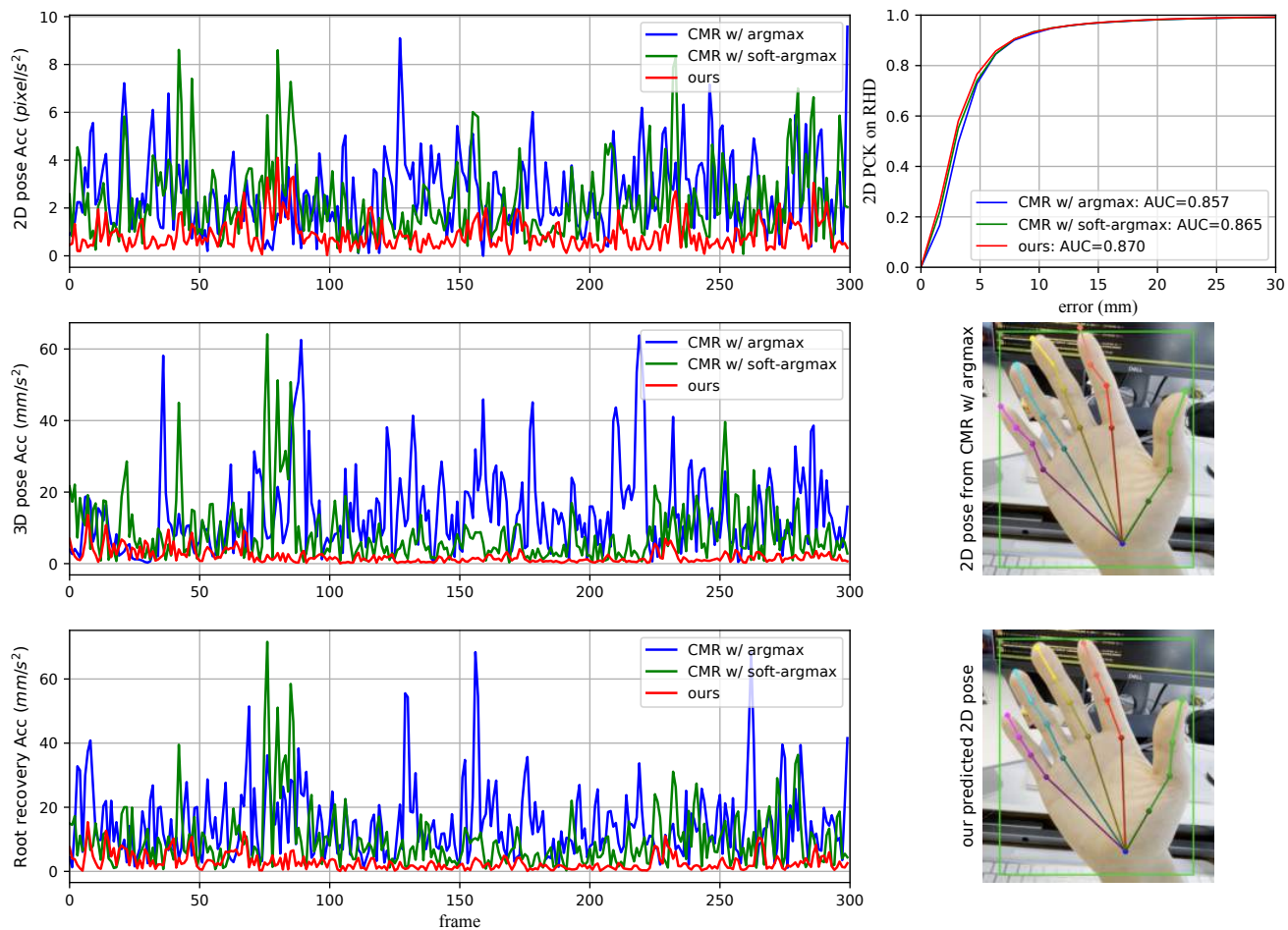


Figure 22. We record a video snippet with a straightforward and static hand pose (see the bottom right corner) to compare the temporal performance and articulated structure.