Consistent 3D Hand Reconstruction in Video via Self-Supervised Learning

Zhigang Tu, *Member, IEEE,* Zhisheng Huang, Yujin Chen, Di Kang, Linchao Bao, *Member, IEEE,* Bisheng Yang, *Senior Member, IEEE*, and Junsong Yuan, *Fellow, IEEE*

Abstract—We present a method for reconstructing accurate and consistent 3D hands from a monocular video. We observe that the detected 2D hand keypoints and the image texture provide important cues about the geometry and texture of the 3D hand, which can reduce or even eliminate the requirement on 3D hand annotation. Accordingly, in this work, we propose S^2HAND , a self-supervised 3D hand reconstruction model, that can jointly estimate pose, shape, texture, and the camera viewpoint from a single RGB input through the supervision of easily accessible 2D detected keypoints. We leverage the continuous hand motion information contained in the unlabeled video data and explore $S^2HAND(V)$, which uses a set of weights shared S^2HAND to process each frame and exploits additional motion, texture, and shape consistency constrains to obtain more accurate hand poses, and more consistent shapes and textures. Experiments on benchmark datasets demonstrate that our self-supervised method produces comparable hand reconstruction performance compared with the recent full-supervised methods in single-frame as input setup, and notably improves the reconstruction accuracy and consistency when using the video training data.

Index Terms—hand pose estimation, 3D hand reconstruction, video analysis, self-supervision

1 INTRODUCTION

TANDS play a central role in the interaction between 2 humans and the environment, from physical contact 3 and grasping to daily communications via hand gesture. Learning 3D hand reconstruction is the preliminary step 5 for many computer vision applications such as augmented 6 reality [1], sign language translation [2], [3], action recog-7 nition [4], [5], and human-computer interaction [6], [7], [8]. However, due to diverse hand configurations and interac-9 tion with the environment, 3D hand reconstruction remains 10 a challenging problem, especially when the task relies on 11 monocular data as input. 12

Compared with multi-view images [9], [10], [11], [12] 13 and depth maps [13], [14], [15], [16], [17], the monocular 14 hand images are more common in practice. In recent years, 15 we have witnessed many efforts in recovering 3D shapes of 16 human hands from single-view RGB images. For example, 17 [18], [19], [20], [21], [22] were proposed to predict 3D hand 18 pose from an RGB image. However, they only represent the 19 3D hand through sparse joints, and ignore the 3D shape 20 information, which are required for some applications such 21 as grasping objects with virtual hands [6]. To better cap-22 ture the surface information of the hand, previous studies 23 predict the triangle mesh either via regressing per-vertex 24 coordinates [23], [24] or by deforming a parametric hand 25 model [25], [26]. Outputting such high-dimensional repre-26

- Z. Tu, Z. Huang, B. Yang are with the State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University, Wuhan 430079, China (email: {tuzhigang, hzs5230, bshyang}@whu.edu.cn). Z. Tu and Z. Huang contributed equally and they are co-first authors.
- Y. Chen is with Technical University of Munich. Work done at Wuhan University. Correspondence to: Y. Chen (email: yujin.chen@tum.de)
- D. Kang and L. Bao are with Tencent AI Lab.
- J. Yuan is with the Computer Science and Engineering Department, University at Buffalo, Buffalo, NY 14228, USA. (email: jsyuan@buffalo.edu).



Fig. 1: Given a collection of unlabeled hand images or videos, we learn a 3D hand reconstruction network in a self-supervised manner. Top: the training uses unlabeled hand images from image collections or video sequences and their corresponding noisy detected 2D keypoints. Bottom: our model outputs accurate hand joints and shapes, as well as vivid hand textures.

sentations from 2D input is challenging for neural networks 27 to learn. As a result, the training process relies heavily on 28 3D hand annotations such as dense hand scans, model-29 fitted parametric hand mesh, or human-annotated 3D joints. 30 Besides, the hand texture is important in some applications, 31 such as vivid hands reconstruction in immersive virtual re-32 ality. But only recent work try to explore parametric texture 33 estimation in a learning-based hand recovery system [27], 34 while most previous work of 3D hand reconstruction do not 35 consider texture modeling [27]. 36

One of our key observations is that the 2D cues in the

hand image are quite informative to reconstruct the 3D hand 38 model in the real world. The 2D hand keypoints contain rich 39 structural information, and the 2D image contains abundant 40 texture and shape information. Both are important for re-41 ducing the use of expensive 3D annotations but have not 42 been fully investigated. Leveraging these cues, we could 43 44 directly use 2D annotations and the input image to learn the geometry and texture representations without relying 45 on 3D annotations [21]. However, it is still labor-consuming 46 to annotate 2D hand keypoints per image. To completely 47 save the manual annotation, we propose to extract 2D hand 48 keypoints as well as geometric representations from the 49 unlabeled hand image to help the shape reconstruction and 50 use the texture information contained in the input image to 51 help the texture modeling. 52

Additionally, video sequences contain rich hand motion 53 and more comprehensive appearance information. Usually, 54 a frame-wise fully-supervised hand reconstruction model 55 does not take these information into serious consideration 56 since 3D annotations already provide a strong supervision. 57 As a result, it is more difficult for a frame-wise model 58 to produce consistent results from video frames compared 59 to sequence-wise models, since no temporal information is 60 utilized. Thereby, we propose to penalize the inconsistency 61 of the output hand reconstructions from consecutive ob-62 servations of the same hand. In this way, motion prior in 63 video is distilled in the frame-wise model to help reconstruct 64 more accurate hand for every single frame. Notably, the 65 constraints on the sequence output are also employed in 66 a self-supervised manner. 67

Driven by the above observations, this work aims to 68 train an accurate 3D hand reconstruction network using 69 only the supervision signals obtained from the input images 70 or video sequences while eliminating manual annotations 71 of the training images. To this end, we use an off-the-72 shelf 2D keypoint detector [28] to generate some noisy 73 2D keypoints, and supervise the hand reconstruction by 74 these noisy detected 2D keypoints and the input image. 75 Although our reconstruction network relies on the pre-76 defined keypoint detector, we call it a self-supervised net-77 work, following the naming convention in the face recon-78 struction literature [29], [30] as only the self-annotation is 79 provided to the training data. Further, we leverage the self-80 supervision signal embedded in the video sequence to help 81 the network produce more accurate and temporally more 82 coherent hand reconstructions. To learn in a self-supervised 83 manner, there are several issues to be addressed. First, how 84 to efficiently use joint-wise 2D keypoints to supervise the ill-85 posed monocular 3D hand reconstruction? Second, how to 86 handle noise in the 2D detection output since our setting 87 is without utilizing any ground truth annotation? Third, 88 is it possible to make use of the continuous information 89 contained in video sequences to encourage smoothness and 90 consistency of reconstructed hands in a frame-wise model? 91

To address the first issue, a model-based autoencoder is 92 learned to estimate 3D joints and shape, where the output 93 3D joints are projected into 2D image space and forced to 94 95 align with the detected keypoints during training. However, if we only align keypoints in image space, invalid hand pose 96 often occurs. This may be caused by an invalid 3D hand 97 configuration which is still compatible with the projected 98

2D keypoints. Furthermore, 2D keypoints cannot reduce the scale ambiguity of the predicted 3D hand. Thus, we propose 100 to learn a series of priors embedded in the model-based 101 hand representations to help the neural network output 102 hand with a reasonable pose and size. 103

To address the second issue, a trainable 2D keypoint 104 estimator and a novel 2D-3D consistency loss are proposed. 105 The 2D keypoint estimator outputs joint-wise 2D keypoints 106 and the 2D-3D consistency loss links the 2D keypoint esti-107 mator and the 3D reconstruction network to make the two 108 mutually beneficial to each other during the training. In 109 addition, we find that the detection accuracy of different 110 samples varies greatly, thus we propose to distinguish each 111 detection item to weigh its supervision strength accordingly. 112

To address the third issue, we decompose the hand mo-113 tion into the joint rotations and ensure smooth rotations of 114 hand joints between frames by conforming to a quaternion-115 based representation. Furthermore, a novel quaternion loss 116 function is proposed to allow all possible rotation speeds. 117 Besides motion consistency, hand appearance is another 118 main concern. A texture and shape (T&S) consistency loss 119 function is introduced to regularize the coherence of the 120 output hand texture and shape. 121

In brief, we present a self-supervised 3D hand reconstruction model $S^{2}HAND$ and its advanced $S^{2}HAND(V)$. The models enable us to train neural networks that can predict 3D pose, shape, texture and camera viewpoint from images without any ground truth annotation of training images, except that we use the outputs from a 2D keypoint detector (see Fig. 1). Notably, $S^{2}HAND(V)$ is able to extract informative supervision from unannotated videos to help learn a better frame-wise model. In order to achieve this, $S^{2}HAND(V)$ inputs the sequential data to multiple weightshared S²HAND models and employs proposed constraints on the sequential output at the training stage.

The advantage of our proposed methods are summarized as follows:

- We present the first self-supervised 3D hand recon-• struction models, which accurately output 3D joints, mesh, and texture from a single image, without using any annotated training data.
- We exploit an additional trainable 2D keypoint estimator to boost the 3D reconstruction through a mutual improvement manner, in which a novel 2D-3D consistency loss is constructed.
- We introduce a hand texture estimation module to learn vivid hand texture via self-supervision.
- We benchmark self-supervised 3D hand reconstruction on some currently challenging datasets, where our self-supervised method achieves comparable performance to previous fully-supervised methods.

This work is an extension of our conference paper [31]. The new contributions include:

- We extend our $S^{2}HAND$ model to the $S^{2}HAND(V)$ • model, which further exploits the self-supervision signals embedded in video sequences. The improvement in accuracy and smoothness is 3.5% and 3.1%, respectively.
- We present a quaternion loss function, which is based 157 on an explored motion-aware joints rotation repre-158

99

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

sentation, to help learn smooth hand motion. Experi ments demonstrate its significant advantage over the
 similar methods in both accuracy and smoothness.

- i sininar metrious in bour accuracy and smoothiess
- We propose a texture and shape consistency regular ization term to encourage coherent shape and texture
 reconstruction.
- We illustrate that utilizing extra in-the-wild unlabeled training data can further boost the performance of our model.

168 2 RELATED WORK

Hand Pose and Shape Estimation. Researchers have devel-169 oped a lot of different methods in hand pose and shape 170 estimation, such as regression-based method [22], [32], [33], 171 [34], [35] and model-based method [36], [37], [38], [39]. Com-172 paring to hand pose which is represented by 3D coordinates 173 of hand joints alone, hand mesh contains more detailed 174 shape information and recently has become the focus in the 175 research community. Several methods utilize the hand mesh 176 topology to directly output 3D mesh vertices. E.g. [40], [41], 177 [42] use the spiral convolution to recover hand mesh and 178 [34], [43], [44] use the graph convolution to output mesh 179 vertices. Although these methods introduce as few priors 180 as possible, they require large amounts of annotated data 181 for training. In this self-supervised work, we make use of 182 the priors contained in the MANO hand model [45], where 183 MANO can map pose and shape parameters to a triangle 184 mesh [26], [46], [47], [48], to reduce reliance on the labeled 185 training data. 186

Because the parametric model contains abundant struc-187 ture priors of human hands, recent works integrate hand 188 model as a differentiable layer in neural networks [25], [26], 189 [46], [48], [49], [50], [51], [52]. Among them, [49], [51], [52] 190 output a set of intermediate estimations, like segmentation 191 mask and 2D keypoints, and then map these representations 192 to the MANO parameters. Different from them, we aim at 193 demonstrating the feasibility of a self-supervised framework 194 using an intuitive autoencoder. We additionally output 2D 195 keypoint estimation from another branch and utilize it only 196 during training to facilitate 3D reconstruction. More gener-197 ally, recent methods [25], [26], [46], [48], [50] directly adopt 198 an autoencoder that couples an image feature encoding 199 stage with a model-based decoding stage. Unlike [25], [26], 200 we focus on hand recovery and do not use any annotation 201 about objects. More importantly, the above methods use 202 3D annotations as supervision, while the proposed method does not rely on any ground truth annotations. 204

3D Hand Pose and Shape Estimation with Limited 205 Supervision. 2D annotation is cheaper than 3D annotation, 206 but it is difficult to deal with the ambiguity of depth and 207 scale. [19] uses a depth map to perform additional weak 208 supervision to strengthen 2D supervision. [21] proposes 209 the biomechanical constraints to help the network output 210 feasible 3D hand configurations. [53] detects 2D hand key-21 points and directly fits a hand model to the 2D detection. 212 [24] gathers a large-scale dataset through an automated 213 214 data collection method similar to [53] and then applies the collected mesh as supervision. In this work, we limit 215 biomechanical feasibility by introducing a set of constraints 216 217 on the skin model instead of only imposing constraints on the skeleton as [21]. In contrast to [19], [24], our method is designed to verify the feasibility of (noisy) 2D supervision and avoids introducing any extra 2.5D or 3D data. 220

Self-supervised 3D Reconstruction. Recently, there are 221 methods that propose to learn 3D geometry from the monoc-222 ular image only. For example, [54] presents an unsupervised 223 approach to learn 3D deformable objects from raw single-224 view images, but they assume the object is perfectly sym-225 metric, which is not the case in the hand reconstruction. 226 [55] removes keypoints from the supervision signals, but 227 it uses ground truth 2D silhouette as supervision and only 228 tackles categories with small intra-class shape differences, 229 such as birds, shoes, and cars. [56] exploits a self-supervised 230 contrastive learning for hand pose estimation, but only the 231 encoder is pretrained in the self-supervised manner. [57] 232 designs a self-supervised module to overcome inconsistency 233 between the 2D and 3D hand pose, but they only consider 234 the sparse joint keypoints. [58] explores a depth-based self-235 supervised 3D hand pose estimation method, but the depth 236 image provides much stronger evidence and supervision 237 than the RGB image. Recently, [29], [30], [59] exploits a self-238 supervised face reconstruction method with the usage of 3D 230 morphable model of face (3DMM) [60] and 2D landmarks 240 detection. Our approach is similar to them, but the hand 241 is non-flat and asymmetrical when compared with the 3D 242 face, and the hand suffers from more severe self-occlusion. 243 These characteristics make this self-supervised hand recon-244 struction task more challenging. 245

Texture Modeling in Hand Recovery. [61], [62] exploit 246 shading and texture information to handle the self-occlusion 247 problem in the hand tracking system. Recently, [27] uses 248 principal component analysis (PCA) to build a parametric 249 texture model of hand from a set of textured scans. In 250 this work, we try to model texture from self-supervised 251 training without introducing extra data, and further inves-252 tigate whether the texture modeling helps with the shape 253 modeling. 254

Motion Learning from Sequence Data for 3D Hand 255 Estimation. To leverage motion information contained in 256 sequence data, several methods have been proposed in 257 hand pose estimation. [25] uses the photometric consistency 258 between neighboring frames of sparsely annotated RGB 259 videos. [63] presents a graph-based method to exploit spatial 260 and temporal relationship for sequence pose estimation. 261 [64] utilizes the temporal information through bidirectional 262 inferences. [36], [65], [66] design a temporal consistency loss 263 for motion smoothness. However, these methods either are 264 specialized for motion generation or only impose a weak 265 regularization for motion smoothness. 266

There exists no approach to capture hand motion dy-267 namics fundamentally, leading to limited benefits can be 268 gained from modeling motion. In this work, we aim to 269 exploit self-supervised information from hand motion dy-270 namics. Unlike most of the previous approaches [67], [68], 271 [69], [70] which adopt recurrent or graph-based network 272 structure to learn hand motion in a sequence-to-sequence 273 manner, we instead use a motion-related loss function to 274 help our frame-wise model converges better and bridges 275 the gap with fully-supervised methods. 276

From the above analysis and comparison, we believe 277 that self-supervised 3D hand reconstruction is feasible and 276



Fig. 2: Overview of the proposed models. The $S^{2}HAND(V)$ on the right learns to reconstruct consistent 3D hands from video sequences without ground truth annotations based on S²HAND. Given an input image, the S²HAND model generates a 3D textured hand with its corresponding multiple 2D representations through a 3D reconstruction network and a 2D keypoints estimator. Effective loss functions and regularization terms are designed for self-supervised network training. Given a video sequence, the $S^{2}HAND(V)$ model produces sequential outputs from several weight-shared $S^{2}HAND$ models with temporal constraints. A quaternion loss and a T&S loss are presented to exploit continuous motion information to promote consistent hand reconstruction. During the inference, only the 3D reconstruction network is utilized and the $S^2HAND(V)$ acts just like a specially trained S²HAND due to weight sharing. The symbols used in this figure can be found in Section 3.2 and Section 3.3.

significant, but to the best of our knowledge, no such idea 279 has been studied in this field. In this work, we fill this gap 280 and propose the first self-supervised 3D hand reconstruc-28 tion model, and prove its effectiveness through extensive 282 experiments. 283

METHODOLOGY 3 284

3.1 Overview 285

Our method enables end-to-end learning of accurate and 286 consistent 3D hand reconstruction from video sequences in a 28 self-supervised manner through $S^{2}HAND(V)$ (Section 3.3), 288 which is based on $S^{2}HAND$ (Section 3.2). The overview is 289 illustrated in Fig. 2. 290

The S²HAND model takes an image as input and gener-291 ates a textured 3D hand represented by pose, shape and tex-292 ture, along with corresponding lighting, camera viewpoint 293 (Section 3.2.1 and 3.2.2) and multiple 2D representations in 294 the image space (Section 3.2.3). Some efficient loss functions 295 and regularization terms (Section 3.2.4) are explored to train 296 the network without using ground truth annotations. The 297 $S^{2}HAND(V)$ model takes video sequences as input and produces consistent sequential outputs from multiple S²HAND 299 models where their weights are shared. A quaternion loss 300 301 (Section 3.3.1) and a T&S consistency loss (Section 3.3.2) are designed to train the network with temporal constraints. We 302 describe the proposed method in detail as below. 303

3.2 Self-supervised Hand Reconstruction from Image 304

Collections 305

The S²HAND model learns self-supervised 3D hand recon-306 struction from image collections via training a 3D hand 307 reconstruction network with the help of a trainable 2D 308 keypoints estimator (See Section 3.2.3). 309

3.2.1 Deep Hand Encoding

Given an image *I* that contains a hand, the 3D hand recon-311 struction network first extracts the feature maps with the 312 EfficientNet-b0 backbone [71], and then transforms them 313 into a geometry semantic code vector x and a texture semantic code vector y. The geometry semantic code vector x parameterizes the hand pose $\theta \in \mathbb{R}^{30}$, shape $\beta \in \mathbb{R}^{10}$, 316 scale $s \in \mathbb{R}^1$, rotation $R \in \mathbb{R}^3$ and translation $T \in \mathbb{R}^3$ in 317 a unified manner: $x = (\theta, \beta, s, R, T)$. The texture semantic 318 code vector y parameterizes the hand texture $C \in \mathbb{R}^{778 \times 3}$ 319 and scene lighting $L \in \mathbb{R}^{11}$ in a unified manner: y = (C, L). 320

3.2.2 Model-based Hand Decoding

Given the geometry semantic code vector x and the texture 322 semantic code vector y, our model-based decoder generates 323 a textured 3D hand model in the camera space. In the fol-324 lowing, we will describe the used hand model and decoding 325 network in detail. 326

Pose and Shape Representation. The hand surface is 327 represented by a manifold triangle mesh $M \equiv (V, F)$ with 328 n = 778 vertices $V = \{v_i \in \mathbb{R}^3 | 1 \le i \le n\}$ and faces 329 F. The faces F indicates the connection of the vertices in 330 the hand surface, where we assume the face topology keeps 331 fixed. Given the mesh topology, a set of k = 21 joints [26] 332 $J = \{j_i \in \mathbb{R}^3 | 1 \le i \le k\}$ (as shown in Fig. 3A) can be 333 directly formulated from the hand mesh. Here, the hand 334 mesh and joints are recovered from the pose vector θ and 335 the shape vector β via MANO, where MANO is a low-336 dimensional parametric model [45]. 337

3D Hand in Camera Space. After representing 3D hand 338 via MANO hand model from pose and shape parameters, 339 the mesh and joints are located in the hand-relative coor-340 dinate systems. To represent the output joints and mesh in 341

314 315

321

the camera coordinate system, we use the estimated scale, rotation and translation to conserve the original hand mesh M_0 and joints J_0 into the final representations in terms of: $M = sM_0R + T$ and $J = sJ_0R + T$.

Texture and Lighting Representation. We use per-face 346 RGB value of 1538 faces to represent the texture of hand 347 $C = \{c_i \in \mathbb{R}^3 | 1 \le i \le n\}$, where c_i yields the RGB values 348 of vertex *i*. In our model, we use a simple ambient light and 349 a directional light to simulate the lighting conditions [72]. 350 The lighting vector L parameterizes ambient light intensity 351 $l^a \in \mathbb{R}^1$, ambient light color $l^a_c \in \mathbb{R}^3$, directional light 352 intensity $l^d \in \mathbb{R}^1$, directional light color $l_c^d \in \mathbb{R}^3$, and di-353 rectional light direction $n^d \in \mathbb{R}^3$ in a unified representation: 354 $L = (l^a, l^a_c, l^d, l^d_c, n^d).$ 355

356 3.2.3 2D Hand Representations

A set of estimated 3D joints within the camera space can 357 be projected into the image space by camera projection. 358 359 Similarly, the output textured model can be formulated into a realistic 2D hand image through a neural renderer. In 360 addition to the 2D keypoints projected from the modelbased 3D joints, we can also estimate the 2D position of 362 each keypoint in the input image. Here, we represent 2D 363 hand with three modes and explore the complementarity 364 among them. 365

Joints Projection. Given a set of 3D joints in camera coordinates J and the intrinsic parameters of the camera, we use perspective camera projection Π to project 3D joints into a set of k = 21 2D joints $J^{pro} = \{j_i^{pro} \in \mathbb{R}^2 | 1 \le i \le k\}$, where j_i^{pro} yields the position of the *i*-th joint in image UV coordinates: $J^{pro} = \Pi(J)$.

Image Formation. A 3D mesh renderer is used to conserve the triangle hand mesh into a 2D image, here we use the implementation¹ of [72]. Given the 3D mesh M, the texture of the mesh C and the lighting L, the neural renderer Δ can generate a silhouette of hand S^{re} and a color image I^{re} : S^{re} , $I^{re} = \Delta(M, C, L)$.

Extra 2D Joint Estimation. Projecting model-based 3D 378 joints into 2D can help the projected 2D keypoints retain 379 the structural information, but at the same time gives 380 up the knowledge of per-joint prior. To address this is-381 sue, we additionally use a 2D keypoint estimator to di-382 rectly estimate a set of k = 21 independent 2D joints 383 $J^{2d} = \{j_i^{2d} \in \mathbb{R}^2 | 1 \le i \le k\}$, where j_i^{2d} indicates the posi-384 tion of the *i*-th joint in image UV coordinates. In our 2D 385 keypoint estimator, a stacked hourglass network [73] along 386 with an integral pose regression [74] is used. Note that the 387 2D hand pose estimation module is optionally deployed in 388 the training period and is not required during the inference. 389

390 3.2.4 Training Objective

Our overall training loss E_{S^2HAND} consists of three parts, i.e. a 3D branch loss E_{3d} , a 2D branch loss E_{2d} and a 2D-3D consistency loss E_{cons} :

$$E_{S^2HAND} = w_{3d}E_{3d} + w_{2d}E_{2d} + w_{cons}E_{cons}$$
(1)

Note, E_{2d} and E_{cons} are optional and only used when the 2D estimator is applied. The constant weights w_{3d} , w_{2d} and

1. https://github.com/daniilidis-group/neural_renderer

To train the model-based 3D hand decoder, we enforce geometric alignment E_{geo} , photometric alignment E_{photo} , and statistical regularization E_{regu} :

$$E_{3d} = w_{geo}E_{geo} + w_{photo}E_{photo} + w_{regu}E_{regu}$$
(2)

Geometric Alignment. We propose a geometric alignment 401 loss E_{qeo} based on the detected 2D keypoints which are 402 obtained through an implementation² of [28]. The detected 403 2D keypoints $L = \{(j_i^{de}, con_i) | 1 \leq i \leq k\}$ allocate each 404 keypoint with a 2D position $j_i^{de} \in \mathbb{R}^2$ and a 1D confi-405 dence $con_i \in [0,1]$. The geometric alignment loss in the 406 2D image space consists of a joint location loss E_{loc} and 407 a bone orientation loss E_{ori} . The joint location loss E_{loc} 408 enforces the projected 2D keypoints J^{pro} to be close to its 409 corresponding 2D detections \hat{J}^{de} , and the bone orientation 410 loss E_{ori} enforces the m = 20 bones of the keypoints in 411 these two sets to be aligned: 412

$$E_{loc} = \frac{1}{k} \sum_{i=1}^{k} con_i \mathcal{L}_{SmoothL1}(j_i^{de}, j_i^{pro})$$
(3)

$$E_{ori} = \frac{1}{m} \sum_{i=1}^{m} conf_i^{bone} \| \nu_i^{de} - \nu_i^{pro} \|_2^2$$
(4)

Here, a SmoothL1 loss [75] is used in Eq. 3 to make the 414 loss term to be more robust to local adjustment since the 415 detected keypoints are not fit well with the MANO key-416 points. In Eq. 4, ν_i^{de} and ν_i^{pro} are the normalized *i*-th bone 417 vector of the detected 2D joints and the projected 2D joints, 418 respectively, and $conf_i^{bone}$ is the product of the confidence 419 of the two detected 2D joints of the *i*-th bone. The overall 420 geometric alignment loss E_{geo} is the weighted sum of E_{loc} 421 and E_{ori} with a weighting factor w_{ori} : 422

$$E_{geo} = E_{loc} + w_{ori} E_{ori} \tag{5}$$

Photometric Consistency. For the image formation, the 423 ideal result is the rendered color image I^{re} matches the 424 foreground hand of the input *I*. To this end, we employ a 425 photometric consistency which has two parts: the pixel loss 426 E_{pixel} is computed by averaging the least absolute deviation 427 (L1) distance for all visible pixels to measure the pixel-wise 428 difference, and the structural similarity (SSIM) loss E_{SSIM} 429 is estimated by evaluating the structural similarity between 430 the two images [76]: 431

$$E_{pixel} = \frac{conf_{sum}}{|S^{re}|} \sum_{(u,v)\in S^{re}} \| I_{u,v} - I_{u,v}^{re} \|_{1}$$
(6)

$$E_{SSIM} = 1 - SSIM(I \odot S^{re}, I^{re}) \tag{7}$$

Here, the rendered silhouette S^{re} is used to get the foreground part of the input image for loss computation. In Eq. 6, we use $conf_{sum}$, which is the sum of the detection confidence of all keypoints, to distinguish different training samples. This is because we think that low-confidence samples correspond to ambiguous texture confidence, e.g., the detection confidence of an occluded hand is usually low. The

 w_{cons} balance the three terms. In the following, we describe these loss terms in detail.

^{2.} https://github.com/Hzzone/pytorch-openpose



Fig. 3: (A)The joint skeleton structure. (B) A sample of bone rotation angles. The five bones $(\overline{01}, \overline{05}, \overline{09}, \overline{013}, \overline{017})$ on the palm are fixed. Each finger has 3 bones, and the relative orientation of each bone from its root bone is represented by azimuth, pitch, and roll.

photometric consistency loss E_{photo} is the weighted sum of 440 441 E_{pixel} and E_{SSIM} with a weighting factor w_{SSIM} :

$$E_{photo} = E_{pixel} + w_{SSIM} E_{SSIM} \tag{8}$$

Statistical Regularization. During training, to make the 443 results plausible, we introduce some regularization terms, 444 including the shape regularization E_{β} , the texture regular-445 ization E_C , the scale regularization E_s , and the 3D joints 446 regularization E_J . The shape regularization term is defined 447 as $E_{\beta} = \| \beta - \overline{\beta} \|$ to encourage the estimated hand model 448 shape β to be close to the average shape $\bar{\beta} = \vec{0} \in \mathbb{R}^{10}$. 449 The texture regularization E_C is used to penalize outlier 450 RGB values. The scale regularization term E_s is used to 451 ensure the output hand has appropriate size, so as to help 452 determining the depth of the output in this monocular 453 3D reconstruction task. To enforce the regularizations on 454 skeleton E_J , we define feasible range for each rotation angle 455 a_i (as shown in Fig. 3B) and penalize those who exceed 456 the feasible threshold. The remaining E_C , E_s and E_J terms 457 follow [31]. 458

The statistical regularization E_{regu} is the weighted sum 459 460 of E_{β} , E_{C} , E_{s} and E_{J} with weighting factors w_{C} , w_{s} and w_{I} : 461

$$E_{regu} = E_{\beta} + w_C E_C + w_s E_s + w_J E_J \tag{9}$$

2D Branch Loss. For the 2D keypoint estimator, we use a joint location loss as in Eq. 3 with replacing the projected 2D 463 joint j_i^{pro} by the estimated 2D joint j_i^{2d} : 464

$$E_{2d} = \frac{1}{k} \sum_{i=1}^{k} con_i \mathcal{L}_{SmoothL1}(j_i^{de}, j_i^{2d})$$
(10)

2D-3D Consistency Loss. Since the outputs of the 2D 465 branch and the 3D branch are intended to represent the 466 same hand in different spaces, they should be consistent 467 when they are transferred to the same domain. Through 468 this consistency, structural information contained in the 3D 469 470 reconstruction network can be introduced into the 2D keypoint estimator, and meanwhile the estimated 2D keypoints 471 can provide joint-wise geometric cues for 3D hand recon-472 struction. To this end, we propose a novel 2D-3D consistency 473

loss to link per projected 2D joint j_i^{pro} with its corresponding 474 estimated 2D joint j_i^{2d} : 475

$$E_{cons} = \frac{1}{k} \sum_{i=1}^{k} \mathcal{L}_{SmoothL1}(j_i^{pro}, j_i^{2d})$$
(11)

3.3 Consistent Self-supervised Hand Reconstruction 476 from Video Sequences

The $S^{2}HAND(V)$ model learns consistent self-supervised 478 3D hand reconstruction from video sequences via training 479 weight-shared S²HAND models with temporal constraints, 480 including a quaternion loss and a T&S consistency loss. 481

3.3.1 Quaternion-based Motion Regularization

We reformulate hand motion in joint rotation perspective. 483 We choose the unit quaternion [77], [78], [79], [80] as our 484 joint rotation representation, which can represent spatial 485 orientations and rotations of elements in a convenient and 486 efficient way. The unit quaternion associated with a spatial 487 rotation is constructed as: 488

$$\mathbf{q} = \left(\cos(\frac{\alpha}{2}), \sin(\frac{\alpha}{2})\vec{\mathbf{u}}\right) \tag{12}$$

where α is the rotation angle and \vec{u} denotes the rotation 489 axis in \mathbb{R}^3 . Notably, **q** can represent both rotation and 490 orientation. 491

Smooth orientation transition q_t between initial q_0 joint 492 orientation and final q1 joint orientation is defined by a 493 unique axis $\vec{\mathbf{v}}$ and corresponding rotation angle γ around 494 the axis. The transition process can be expressed as follows: 495

$$\mathbf{q}_t = (\mathbf{q}_1 \mathbf{q}_0^{-1})^{\phi(t)} \mathbf{q}_0 = \left(\cos(\frac{\gamma}{2}), \sin(\frac{\gamma}{2}) \vec{\mathbf{v}}\right)^{\phi(t)} \mathbf{q}_0 \qquad (13)$$

where \mathbf{q}_0^{-1} represents the inverse of \mathbf{q}_0 , and the product 496 operation here is the Hamilton product. The $\phi(t)$ denotes a 497 monotonically non-decreasing function, ranging from 0 to 1 498 and controlls the orientation transition from q_0 to q_1 . When 499 $\phi(t)$ equals 0 or 1, \mathbf{q}_t will equal to \mathbf{q}_0 and \mathbf{q}_1 respectively. 500 In order to reduce the computational cost brought by the 501 Hamilton product, we further rewrite Eq. 13 as a linear 502 combination of the two quaternions \mathbf{q}_0 and \mathbf{q}_1 : 503

$$\mathbf{q}_{t} = (\mathbf{q}_{1}\mathbf{q}_{0}^{-1})^{\phi(t)}\mathbf{q}_{0} = Norm\left[\mu(t)\mathbf{q}_{0} + \varepsilon(t)\mathbf{q}_{1}\right]$$
(14)

where $Norm[\cdot]$ denotes normalization to ensure the result 504 is a unit quaternion. $\mu(t)$ and $\varepsilon(t)$ are time-dependent 505 coefficients which are determined by $\phi(t)$. One instance of 506 Eq. 14 is Slerp [79], which is a widely used linear quaternion 507 interpolation method with constant rotation speed, assum-508 ing $\phi(t) = t$: 509

$$\mathbf{q}_t = (\mathbf{q}_1 \mathbf{q}_0^{-1})^t \mathbf{q}_0 = \frac{\sin((1-t)\eta)}{\sin(\eta)} \mathbf{q}_0 + \frac{\sin(t\eta)}{\sin(\eta)} \mathbf{q}_1 \qquad (15)$$

where η is the included angle between \mathbf{q}_0 and \mathbf{q}_1 as two 510 vectors, which can be computed by: 511

$$\eta = \frac{\mathbf{q}_1 \cdot \mathbf{q}_0}{\|\mathbf{q}_1\| \, \|\mathbf{q}_0\|} \tag{16}$$

where \cdot denotes inner product of two vectors and $\|\cdot\|$ is the 512 magnitude of a vector. 513

Instead of generating interpolated poses as psuedo-514 labels with one specific $\phi(t)$ for supervision, we propose 515

477



Fig. 4: Comparison between our quaternion loss and Slerp. The circle represents a 2D projective plane of 4D unit quaternion sphere. The red arch denotes the set of quaternion that satisfies Eq. 13, which ensures smooth orientation transition. The equation in each circle represents the corresponding prior. The remaining symbols can be found in Section 3.3. As can be seen, both Slerp and quaternion loss has the prior to make sure Eq. 13 is satisfied. However Slerp has an additional prior $\phi(t) = t$, while our Quaternion loss covers all possible $\phi(t)$, which allows smooth orientation transition at all possible speed.

a quaternion loss function to cover all possible joint rotationspeeds as following:

$$E_{quat} = \left\| \sum_{i=1}^{n-1} \Psi(H_i, H_{i+1}) - \Psi(H_1, H_n) \right\|$$
(17)

where Ψ is the function to compute the rotation angle γ between two quaternions, and H_i denotes the output hand pose represented in quaternion of frame *i*. In practice, H_i is the concatenation of i-th pose vector θ_i and i-th rotation R_i to cover all 21 hand joints:

$$H_i = Quaternion(Concatenate[\theta_i, R_i])$$
(18)

where *Quaternion* denotes the transformation from representation of MANO outputs to quaternion representation and *Concatenate* denotes the concatenation operation. The comparison between the proposed quaternion loss and Slerp is illustrated in Fig. 4.

To understand the proposed quaternion loss, two points are important. One is that the quaternion interpolation is essentially finding a rotation curve through a fixed rotation axis between two poses, as suggested by Eq. 13. The other is that the rotation angle γ in the quaternion space relates to included angle η in vector space, as indicated by Eq. 14. Specifically:

$$\cos(\frac{\gamma}{2}) = \cos(\eta) \tag{19}$$

which provides an efficient way to compute γ and is derived from:

$$\Re\left(\mathbf{q}_{1}\mathbf{q}_{0}^{-1}\right) = \Re\left(\cos(\frac{\gamma}{2}), \sin(\frac{\gamma}{2})\vec{\mathbf{v}}\right) = \|\mathbf{q}_{1}\| \|\mathbf{q}_{0}\|\cos(\eta)$$
(20)

where \Re represents the real part of a quaternion, \mathbf{q}_0 , \mathbf{q}_1 , γ and $\vec{\mathbf{v}}$ are the same as before. $\|\cdot\|$ is the magnitude of a vector. η denotes the included angle between \mathbf{q}_0 and \mathbf{q}_1 as two vectors. Eq. 20 can be deduced by comparing the inner product and the Hamilton product of the \mathbf{q}_0 and \mathbf{q}_1 .

3.3.2 T&S Consistency Regularization

We introduce a regularization term on texture and shape to consider consistency of hand appearance in videos. Since texture is coupled with light, our T&S loss is formulated as: 543

$$E_{T\&S} = \sum_{i=1}^{n} \left\| C_i^L - \overline{C^L} \right\| + \sum_{i=1}^{n} \left\| \beta_i - \overline{\beta} \right\|$$
(21)

where C_i^L and β_i are the i-th lighted texture and shape of the sequential reconstruction output, $\overline{C_i^L}$ and $\overline{\beta_i}$ are the corresponding average of the sequential output. The lighted texture C_i^L is computed following [72]:

$$C_{i}^{L} = (l^{a}l_{c}^{a} + (n^{d} \cdot n_{i})l^{d}l_{c}^{d})C_{i}$$
(22)

where n_i is the normal direction of C_i in canonical zero pose [45] and the rest are defined in Section 3.2.1.

A low standard deviation of sequential hand appearance reconstruction from video sequences is promoted by this loss function.

3.3.3 Training Objective

Our overall training loss $E_{S^2HAND(V)}$ consists of three parts, including a S²HAND loss E_{S^2HAND} , a quaternion loss E_{quat} and a T&S loss $E_{T\&S}$:

$$E_{S^2HAND(V)} = E_{S^2HAND} + w_{quat}E_{quat} + w_{ts}E_{T\&S}$$
(23)

where E_{S^2HAND} is the same as that in Section 3.2.4. For E_{quat} and $E_{T\&S}$, please refer to Section 3.3.1 and Section 3.3.2 respectively. The constant weights w_{quat} and w_{ts} are used to balance the three terms.

4 EXPERIMENTS

4.1 Datasets

We evaluate the proposed methods on three datasets. Two of them (FreiHAND and HO-3D) are challenging realistic datasets, aiming for assessing 3D joints and 3D meshes with hand-object interaction. The results are reported through the online submission systems ^{3,4}. The remaining one (STB) is a hand-only video dataset. Besides, we adopt another dataset (YT 3D) to provide in-the-wild data.

The FreiHAND dataset [48] is a large-scale real-world dataset, which contains 32,560 training samples and 3,960 test samples. For each training sample, one real RGB image and extra three images with different synthetic backgrounds are provided. Part of the sample is a hand grabbing an object, but it does not provide any annotations for the foreground object, which poses additional challenges.

The HO-3D dataset [81] collects color images of a hand 579 interacting with an object. The dataset is made of 68 se-580 quences, totaling 77,558 frames of 10 users manipulating 581 one among 10 different objects. The training set contains 582 66,034 images and the test set contains 11,524 images. The 583 objects in this dataset are larger than that in FreiHAND, 584 thus resulting in larger occlusions to hands. We use this 585 dataset in two cases. In the case of self-supervised hand 586 reconstruction from image collections with S²HAND, we do 587 not use the sequence information provided by HO-3D and 588

3. https://competitions.codalab.org/competitions/21238

4. https://competitions.codalab.org/competitions/22485

45

542

552 553 554

555

563

564

550

mix all sequences as a image collection. In the case of self supervised hand reconstruction from video sequences with
 S²HAND(V), we make use of the sequence information
 and compose the training batch accordingly. Details can be
 found in Section 4.3.

The STB dataset [82] is a hand-only dataset, which contains 12 sequences with 18000 frames in total. RGB images along with depth-images, 2D and 3D joint annotations are provided. We follow the splits in [82], using 10 sequences for training and 2 sequences for evaluation. We select this dataset to validate the proposed methods in the hand-only scenario.

The YT 3D dataset [40] contains 116 in-the-wild videos, 601 which is comprised of 102 train videos, 7 validation videos 602 and 7 test videos, along with 47125, 1525 and 1525 hand 603 annotations. We only use this dataset as extra in-the-wild 604 training data in Section 4.5.3 without any annotations. Since 605 606 the 102 train videos are edited with lots of cutaway and only part of videos are accessible due to copyright issues, we 607 preprocess this dataset by filtering out unavailable videos 608 and discontinuous hand motion frames according to the 609 detected 2D keypoints, yielding 34 train videos containing 610 21628 frames with detected 2D keypoints. 611

612 4.2 Evaluation Metrics

We evaluate 3D hand reconstruction by evaluating 3D joints 613 and 3D meshes. For 3D joints, we report the mean per joint 614 position error (MPJPE) in the Euclidean space for all joints 615 on all test frames in *cm* and the **area under the curve** (AUC) 616 of the PCK AUC_J. Here, the PCK refers to the percentage 617 of correct keypoints. For 3D meshes, we report the mean 618 per vertex position error (MPVPE) in the Euclidean space 619 for all joints on all test frames in cm and the AUC of 620 the percentage of correct vertex AUC_V . We also compare 621 the F-score [83] which is the harmonic mean of recall and 622 precision for a given distance threshold. We report distance 623 threshold at 5mm and 15mm and report F-score of mesh 624 vertices at 5mm and 15mm by F_5 and F_{15} . Following the 625 previous works [48], [81], we compare aligned prediction 626 results with Procrustes alignment, and all 3D results are 627 evaluated by the online evaluation system on FreiHAND 628 and HO-3D. For 2D joints, we report the MPJPE in pixel 629 and the curve plot of fraction of joints within distance. 630 For smooth hand reconstruction, we report the acceleration 631 error (ACC-ERR) and the acceleration(ACC) which are first 632 proposed in [84]. ACC-ERR measures average difference 633 between ground truth acceleration and the acceleration of 634 the predicted 3D joints in mm/s^2 while ACC calculates 635 mean acceleration of the predicted 3D joints in mm/s^2 . 636 Generally, lower ACC-ERR and ACC indicate smoother 637 sequence predictions. For shape and texture consistency 638 in sequence predictions, we report corresponding standard 639 deviations (S.D.), in which low deviation means coherent 640 and consistent sequence predictions. Specifically, we com-641 pute texture S.D. and shape S.D., which are the average 642 of per dimensional S.D. of the lighted textures and shape 643 parameters in sequence predictions respectively. 644

645 4.3 Implementation Details

Pytorch [85] is used for implementation. For the 3D reconstruction network, the EfficientNet-b0 [71] is pre-trained on



Fig. 5: Qualitative comparison to OpenPose [28] and MANO-CNN on the FreiHAND testing set. For OpenPose, we visualize the detected 2D keypoints. For our method and MANO-CNN, we visualize both the projected 2D keypoints and 3D mesh.

the ImageNet dataset. The 2D keypoint estimator along with 648 the 2D-3D consistency loss is optionally used. If we train 649 the whole network with the 2D keypoint estimator, a stage-650 wise training scheme is used. We train the 2D keypoint 651 estimator and 3D reconstruction network by 90 epochs 652 separately, where E_{3d} and E_{2d} are used, respectively. The 653 initial learning rate is 10^{-3} and reduced by a factor of 2 654 after every 30 epochs. 655

When training $S^{2}HAND$ with image collections, we 656 finetune the whole network with E_{S^2HAND} by 60 epochs 657 with the learning rate initialized to 2.5×10^{-4} , and reduced 658 by a factor of 3 after every 20 epochs. We use Adam [86] 659 to optimize the network weights with a batch size of 64. 660 We train our model on two NVIDIA Tesla V100 GPUs, 661 which takes around 36 hours for training on FreiHAND. 662 Otherwise, when training $S^{2}HAND(V)$ with $E_{S^{2}HAND(V)}$, 663 the input remains 4D (B,C,H,W), but the sampling strategy 664 is different. Specifically, for a training batch, we first ran-665 domly sample m sequences from the training sequences, 666 then randomly sample n frames in each of the sampled 667 sequences, finally combine these frames to compose a batch. 668 This results in a batch size of *mn*. In our experiments, we set 669 m equal to 64//n, where 64 follows the batch size in training 670 $S^{2}HAND$ and // represents floor division. We find *n* equals 671 3 (see Section 4.5.2 for the ablation study) gets the best 672 performance, thus we use it as the default setting. Notice 673 that all the sampling procedures take place during training. 674 For the stability of this sampling strategy, please refer to 675 Section 2 of the supplementary materials. The learning rate, 676 the reducing schedule, and the Adam optimizer are set the 677 same as before. For the weighting factors, we set $w_{3d} = 1$, 678 $w_{2d} = 0.001, w_{cons} = 0.0002, w_{geo} = 0.001, w_{photo} = 0.005,$ 679 $w_{quat} = 0.05, w_{ts} = 0.01, w_{regu} = 0.01, w_{ori} = 100,$ 680 $w_{SSIM} = 0.2, w_{C} = 0.5, w_{s} = 10$ and $w_{J} = 10$. Here 681 the weighting factors need to regard different fundamental 682 units of the losses, thus the magnitudes do not strictly 683 imply the importance. For how to select these weights and 684 the weights sensitivity, please refer to the supplementary 685 materials. 686

TABLE 1: Comparison of main results on the FreiHAND testing set. The performance of our self-supervised method S^2 HAND is comparable to the recent fully-supervised and weakly-supervised methods. [21]* also uses the synthetic training data with 3D supervision. Note that FreiHAND is not presented with video sequences, which disables learning of S^2 HAND(V)

upervision. Note that frem fAND is not presented with video sequences, which disables fearning of 6 minub (v).							
Supervision	Method	$AUC_{J}\uparrow$	MPJPE↓	$AUC_V\uparrow$	MPVPE↓	$F_5\uparrow$	$F_{15}\uparrow$
	[48](2019)	0.35	3.50	0.74	1.32	0.43	0.90
2D	[26](2019)	0.74	1.33	0.74	1.33	0.43	0.91
3D	[46](2019)	0.78	1.10	0.78	1.09	0.52	0.93
	[27](2020)	0.78	<u>1.11</u>	0.78	<u>1.10</u>	<u>0.51</u>	0.93
2D	[21](2020)*	0.78	1.13	-	-	-	-
-	$S^{2}HAND$	0.77	1.18	0.77	1.19	0.48	0.92

TABLE 2: Comparison of main results on the HO-3D testing set. Without using any object information and hand annotation, our $S^{2}HAND$ model performs comparable with the recent fully-supervised methods [25]. Further with the temporal constraints, our $S^{2}HAND(V)$ even surpasses [25].

Supervision	Method	$AUC_{J}\uparrow$	MPJPE↓	$AUC_V\uparrow$	MPVPE↓	$F_5\uparrow$	$F_{15}\uparrow$
	[26](2019)	-	-	-	1.30	0.42	0.90
3D	[81](2020)	-	-	-	1.06	0.51	0.94
	[25](2020)	<u>0.773</u>	<u>1.11</u>	0.773	1.14	0.43	<u>0.93</u>
	$S^{2}HAND$	0.773	1.14	0.777	1.12	0.45	0.93
-	$S^{2}HAND(V)$	0.780	1.10	0.781	<u>1.09</u>	<u>0.46</u>	0.94

TABLE 3: Ablation studies on different losses used in our method on the FreiHAND testing set. Refer to Section 4.5.1 for details.

		Loss	es		MPIPE	MPVPF	AUCT	$\Delta UC_{12} \uparrow$	F- ↑	F. *
E_{loc}	E_{regu}	E_{ori}	E_{2d} , E_{con}	E_{photo}	1 1 11 J1 L4	IVII VI L _V	nooj	1000	12	1 15
\checkmark					1.97	2.31	0.611	0.545	0.257	0.763
\checkmark	\checkmark				1.54	1.58	0.696	0.687	0.387	0.852
\checkmark	\checkmark	\checkmark			1.24	1.26	0.754	0.750	0.457	0.903
\checkmark	\checkmark	\checkmark		\checkmark	1.22	1.24	0.759	0.754	0.468	0.909
\checkmark	\checkmark	\checkmark	\checkmark		1.19	1.20	0.764	0.763	0.479	0.915
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	1.18	1.19	0.766	0.765	0.483	0.917

687 4.4 Comparison with State-of-the-art Methods

We give comparison on FreiHAND with four recent model-688 based fully-supervised methods [26], [27], [46], [48] and a 689 state-of-the-art weakly-supervised method [21] in Table 1. 690 Note that [87] is not included here since it designs an 69 advanced "image-to-lixel" prediction instead of directly 692 regressing MANO parameters. Our approach S²HAND 693 focuses on providing a self-supervised framework with 694 lightweight components, where the hand regression scheme 695 is still affected by highly non-linear mapping. Therefore, we make a fairer comparison with popular model-based 697 methods [26], [27], [46], [48] to demonstrate the performance 698 of this self-supervised approach. Without using any annota-699 tion, our approach $S^{2}HAND$ outperforms [26], [48] on all 700 evaluation metrics and achieves comparable performance 701 to [27], [46]. [21] only outputs 3D pose, and its pose per-702 formance is slightly better than our results on FreiHAND 703 test set but with much more training data used including 704 RHD dataset [22] (with 40,000+ synthetic images and 3D 705 annotations) as well as 2D ground truth annotation of the 706 FreiHAND. 70

In the hand-object interaction scenario, we compare 708 with three recent fully-supervised methods on HO-3D in 709 Table 2. Compared to the hand branch of [26], both of our 710 self-supervised models S²HAND and S²HAND(V) show 711 higher mesh reconstruction performance where we get a 712 14% and 16% reduction in MPVPE respectively. Compared 713 with [25], which is a fully-supervised joint hand-object pose 714 estimation method, our S²HAND obtains comparable joints 715 and shape estimation results, while our $S^{2}HAND(V)$ even 716 surpasses [25]. [81] gets slightly better shape estimation 717 results than ours, probably due to the fact that they first 718 estimate 2D keypoint positions using heatmaps and then fit 719 MANO model to these keypoints. Since [21] utilizes a totally 720



Fig. 6: Comparisons between our proposed self-supervised (S.S.) methods and other SOTA fully-supervised (F.S.) methods on the STB dataset (hand-only scenario).

different version of HO-3D that is published with HANDS 2019 Challenge⁵, we do not compare with it on HO-3D (not shown in Table 2) as we do on FreiHAND.

In the hand-only scenario, we compare with the fully-724 supervised methods [20], [22], [53], [88], [89] on STB in Fig. 6. 725 Some state-of-the-art fully-supervised methods (e.g. [50] 726 with 0.995 AUC and [90] with 0.997 AUC) are not included 727 in Fig. 6, because their PCK values are not provided and 728 they have similar performance with [20]. The proposed self-729 supervised method S²HAND outperforms some previous 730 fully-supervised methods [22], [53], [88], and $S^{2}HAND(V)$ 731 further improves the performance of S²HAND. On the other 732 side, there is a certain gap between ours and the recent 733 fully-supervised methods [20], [89] since they have direct 734

5. https://sites.google.com/view/hands2019/challenge

721

722

TABLE 4: Comparison of the accuracy of different motion-related constraints on the HO-3D testing set. Quaternion loss shows its effectiveness over the similar loss functions on modeling smoothness.

Method		$AUC_{J}\uparrow$	MPJPE↓	$AUC_V\uparrow$	MPVPE↓	$F_5\uparrow$	$F_{15}\uparrow$
S ² HAND		0.773	1.14	0.777	1.12	0.45	0.93
$S^{2}HAND(V)$ (w/ Temporal Loss	5 [65])	0.771	1.15	0.773	1.14	0.44	0.93
$S^{2}HAND(V)$ (w/ Smooth Loss	[36])	0.774	1.13	0.776	1.12	0.45	0.93
$S^{2}HAND(V)$ (w/ Quaternion L	loss)	0.779	1.10	0.781	1.09	0.46	0.94
$S^{2}HAND(V)$ (w/ Quaternion Loss, T	Г&S Loss)	0.780	1.10	0.781	1.09	0.46	0.94

TABLE 5: Comparison of the smootheness performance of different motion-related constraints on the HO-3D dataset. Quaternion loss gives the smoothest predictions and is highly in line with ACC and ACC-ERR.

Method		Train s	Test set		
Method	ACC-ERR↓	ACC↓	Quaternion Loss↓	ACC↓	Quaternion Loss↓
S ² HAND	2.85	2.29	0.014	3.68	0.020
$S^{2}HAND(V)$ (w/ Temporal Loss [65])	2.57	1.96	0.011	2.87	0.015
$S^{2}HAND(V)$ (w/ Smooth Loss [36])	2.89	2.29	0.014	2.89	0.015
$S^{2}HAND(V)$ (w/ Quaternion Loss)	2.23	1.59	0.008	2.81	0.014



Maximum allowed distance to GT (Pixel) Fig. 7: A comparison of 2D keypoint sets used or outputted at the training stage on FreiHAND. The fraction of frames within the maximum joint distance is plotted. Refer to Section 4.5.1 for details.

735 3D supervision.

In Fig. 5, we show 2D keypoint detection from Open-Pose [28] and our S²HAND results of difficult samples. We also compare the reconstruction results with MANO-CNN, which directly estimates MANO parameters with a CNN [48], but we modify its backbone to be the same as ours. Our results are more accurate and additionally with texture.

742 4.5 Self-comparison

In Section 4.5.1 and Section 4.5.2, we conduct extensive selfcomparisons to verify the effectiveness of each component
of our models S²HAND and S²HAND(V). In Section 4.5.3,
we explore the effect of absorbing extra in-the-wild video
sequences with S²HAND(V).

748 4.5.1 Hand Reconstruction from Image Collections

For self-supervised hand reconstruction from image collection with S²HAND, we conduct ablation studies on FreiHAND, since it is a widely used challenging dataset for
hand pose and shape estimation from single RGB images.

First, we give evaluation results on FreiHAND of settings with different components along with corresponding loss terms used in the network in Table 3. The baseline only uses the 3D branch with E_{loc} and E_{regu} , then we add E_{ori} which helps the MPJPE and MPVPE decrease by 19.5%. After adding the 2D branch with E_{2d} and the 2D-3D

TABLE 6: Comparison of self-supervised results and weaklysupervised results. Refer to Section 4.5.1 for details.

Dataset	Method	$AUC_{J}\uparrow$	$AUC_V\uparrow$	$F_5\uparrow$	$F_{15}\uparrow$		
FroiHAND	WSL	0.730	0.725	0.42	0.89		
FIEITIAND	SSL	0.766	0.765	0.48	0.92		
UO 2D	WSL	0.765	0.769	0.44	0.93		
110-3D	SSL	0.773	0.777	0.45	0.93		

consistency constrain E_{cons} , the MPJPE and MPVPE further reduce by 4%. The E_{photo} slightly improves the pose and shape estimation results.

Then, we make comparison of different 2D keypoint 762 sets. In our approach, there are three sets of 2D keypoints, 763 including detected keypoints J^{de} , estimated 2D keypoints 764 J^{2d} , and output projected keypoints J^{pro} , where J^{de} is 765 used as supervision terms while J^{2d} and J^{pro} are output 766 items. In our setting, we use multiple 2D representations 767 to boost the final 3D hand reconstruction, so we do not 768 advocate the novelty of 2D hand estimation, but compare 769 2D accuracy in the training set to demonstrate the effect 770 of learning from noisy supervision and the benefits of the 771 proposed 2D-3D consistency. Although we use OpenPose 772 outputs as the keypoint supervision source (see *OpenPose* 773 in Fig. 7), we get lower overall 2D MPJPE when we pre-774 train the 2D and 3D branches separately (see Predicted w/o 775 2D-3D and Projected w/o 2D-3D in Fig. 7). After finetuning 776 these two branches with 2D-3D consistency, we find both 777 of them gain additional benefits. After the finetuning, the 778 2D branch (*Predicted w/ 2D-3D*) gains 5.4% reduction in 2D 779 MPJPE and the 3D branch (*Projected wl 2D-3D*) gains 9.3% 780 reduction in 2D MPJPE. From the curves, we can see that 781 2D keypoint estimation (including OpenPose and our 2D 782 branch) gets higher accuracy in small distance thresholds 783 while the regression-based methods (*Projected wlo 2D-3D*) 784 get higher accuracy with larger distance threshold. From 785 the curves, the proposed 2D-3D consistency can improve 786 the 3D branch in all distance thresholds, which verifies the 787 rationality of our network design. 788

Last, we compare the weak-supervised (WSL) scheme using ground truth annotations with our self-supervised (SSL) approach to investigate the ability of our method to handle noisy supervision sources. Both settings use the same network structure and implementation, and WSL uses the ground truth 2D keypoint annotations whose keypoint confidences are set to be the same.

As shown in Table 6, our SSL approach has better performance than WSL settings on both datasets. We think this is

759

760



Fig. 8: Qualitative comparison of different motion-related constraints on the HO-3D testing set. Our $S^{2}HAND(V)$ with the quaternion loss achieves the best qualitative results.

because the detection confidence information is embedded 798 into the proposed loss functions, which helps the network 799 discriminate different accuracy in the noisy samples. In 800 addition, we find that the SSL method outperforms the WSL 801 method in a smaller amplitude on HO-3D (by 1.0%) than 802 that on FreiHAND (by 4.9%). We think this is because the 803 HO-3D contains more occluded hands, resulting in poor 804 2D detection results. Therefore, we conclude that noisy 805 2D keypoints can supervise shape learning for the hand 806 reconstruction task, while the quality of the unlabeled image 807 also has a certain impact. 808

4.5.2 Consistent Hand Reconstruction from Video Sequences

For consistent self-supervised hand reconstruction from video sequences with S²HAND(V), we conduct ablation studies on HO-3D since it is a widely used challenging dataset that presents hand-object interaction images with sequence information.

We first study the accuracy performance of quaternion 816 loss in comparison with other commonly used motion-817 related constraints modeling smoothness in sequence out-818 puts. The compared constraints include temporal loss from 819 [36] closing neighboring poses as much as possible and 820 smooth loss from [65] limiting neighboring pose variation 821 with a threshold. Quantitative and qualitative results are 822 presented in Table 4 and Fig. 8. Compared to S²HAND, 823 $S^{2}HAND(V)$ with quaternion loss improves single frame 824 prediction by reducing 3.5% in MPVPE, while $S^{2}HAND(V)$ 825 with smooth loss from [36] only gets a reduction of less than 826 1% in MPVPE and $S^{2}HAND(V)$ with temporal loss from [65] even degenerates the accuracy. We think this is because 828 weak supervision is prone to optimize models in the wrong 829 direction and suffers from being too sketchy and limited 830

TABLE 7: Comparison of different configurations of the quaternion loss on the HO-3D testing set.

Config	$AUC_{J}\uparrow$	MPJPE↓	$AUC_V\uparrow$	MPVPE↓
interv=1, n=3	0.778	1.11	0.780	1.10
interv=3, n=3	0.780	1.10	0.782	1.09
interv=6, n=3	0.775	1.13	0.776	1.12
interv=3, n=6	0.774	1.13	0.777	1.12

under self-supervised settings. The temporal loss from [65] 831 is beneficial when the 3D annotation is available but may 832 collapse the model by making the network insensitive to 833 the high-frequency details in absence of strong supervision 834 signals. The smooth loss from [36] introduces threshold, but 835 in exchange enlarges the solution space and vanishes when 836 the threshold is exceeded. In contrast, quaternion loss nar-837 rows the solution space based on hand structures and hand 838 motion dynamics and provides more significant supervision 839 signals, which proves its effectiveness over other similar 840 constraints in accuracy. 841

We next explore our quaternion loss with different con-842 figurations in terms of the actual frame interval (interv) of 843 sampled frames and the number of frames (See n in Eq. 17). 844 The results are shown in Table 7. A medium interval value 845 (interv = 3 in Table 7) achieves the best performance. We 846 attribute this to two reasons. On one hand, large interval 847 may witness changes of the joint rotation axis in sampled 848 frames, which contradicts the fixed axis prior in quaternion 849 interpolation (See Section 3.3.1). On the other hand, small 850 interval only witnesses small rotation angles, which limits 851 the effect of quaternion loss. Also, interval is related to 852 the motion speed. In fact, a slow hand motion with larger 853 interval would be equivalent to a fast hand motion with 854 smaller interval. Though the quaternion loss is designed to 855 handle hand motion speed variation, an optimum interval 856 still exists based on the overall motion speed. For number 857 of frames n, increasing it causes model accuracy to drop. 858 We believe this is because optimizing multiple unconfident 859



Fig. 9: Qualitative demonstration of the effectiveness of the T&S consistency loss. For each frame, we show output keypoints (left), output 3D reconstruction (middle), and both sides of the output textures (with lighting) in flat hands (top-/bottom-right). For each sequence, we show results without T&S loss on the top row and with T&S loss on the bottom row. T&S loss significantly improves the output appearance consistency in sequence predictions.



Fig. 10: Boxplots of shape parameters in sequence predictions on the HO-3D testing set. SM1, AP13, and AP10 are sequences from HO-3D testing set. T&S loss reduces S.D. across all 10 dimensions of the shape parameters.

predictions at the same time puts an extra burden on 860 the gradient-descent-based optimizer and destabilizes the 861 learning procedure. From above, we conclude that the best 862 configuration of quaternion loss depends on the frame rates 863 of input video sequences and the confidence of the output. 864 Adjusting the configuration of quaternion loss dynamically 865 according to the input and the output can be a promising 866 direction for future work. 86

We then compare smoothness performance of quater-868 nion loss with others by concatenating their single frame 869 predictions to corresponding sequences predictions. The 870 results are reported on both train set and test set in Table 5. 87 Note that we only report ACC-ERR on train set since 3D 872 ground truth is required to calculate ACC-ERR. As shown 873 in Table 5, all smoothness related prior improve the smooth-874 ness of the sequence outputs, and our proposed quaternion 875 loss achieves the best performance in all evaluation matri-876 877 ces. [65] comes second in terms of smoothness, but shows

the worst performance on accuracy in Table 4. We think 878 better smoothness does not imply higher accuracy, where 879 a trade-off between smoothness and accuracy is shown 880 in some cases. Our quaternion loss differently does best 881 on both accuracy and smoothness, which again proves its 882 superiority over similar methods. In addition, we report 883 the average quaternion loss of the concatenated sequences 884 predictions. We find that our quaternion loss is highly 885 in line with ACC and ACC-ERR, which is encouraging 886 since they are calculated from completely different angles. 887 ACC-ERR and ACC regard no hand structure and hand 888 motion characteristic, rely solely on mechanics in terms of 889 acceleration, while quaternion loss does the opposite. Thus, 890 we think our proposed quaternion loss not only proves its 891 advantages over other loss functions in smoothness, but 892 is capable of being the metric measuring smoothness of 893 sequence predictions as well. 894

Finally, we inspect the effect of regularizing outputs of



Fig. 11: Qualitative results of $S^{2}HAND(V)$ with extra in-the-wild data. The first three rows show results from HO-3D and the last three rows show results from YT 3D.

TABLE 8: Results of absorbing extra in-the-wild data from YT 3D on the HO-3D testing set. * represents changing camera model from perspective to the orthogonal model, which enables to learn with in-the-wild data without camera information.

Method	AUC _J ↑	MPJPE↓	$AUC_V\uparrow$	MPVPE↓	$F_5\uparrow$	$F_{15}\uparrow$
S ² HAND	0.773	1.14	0.777	1.12	0.45	0.93
$S^{2}HAND^{*}$	0.770	1.15	0.774	1.13	0.45	0.93
$S^{2}HAND(V)$ * (w/ Quaternion Loss)	0.778	1.11	0.780	1.10	0.45	0.94
$S^{2}HAND(V)^{*}$ (w/ Quaternion Loss) + YT 3D data	0.782	1.09	0.783	1.09	0.46	0.94

TABLE 9: Results of hand appearance consistency of our methods on the HO-3D testing set. Texture S.D. is computed using lighted textures defined in Eq. 22.

Method	Texture S.D.	Shape S.D.
S ² HAND	0.033	0.013
$S^{2}HAND(V)$ (w/ Quaternion Loss)	0.042	0.016
S ² HAND(V) (w/ Quaternion Loss, T&S Loss)	0.016	0.012

hand shape and texture in sequence predictions with the 896 T&S consistency loss. The results are presented in Table 4, 897 Table 9, Fig. 9 and Fig. 10. Though the proposed T&S consis-898 899 tency loss does not further improve pose accuracy as shown 900 in Table 4, it significantly improves the hand appearance consistency in sequence predictions. Quantitatively, with the 901 T&S loss, the shape S.D. and texture S.D. drops 62% and 902 25% respectively (see Table 9). Qualitatively, $S^{2}HAND(V)$ 903 with the T&S loss gives more consistent reconstructions 904 of different frames from the same sequence than without 905 T&S as shown in Fig. 9. We also provide boxplots of 906 shape parameters to show dimensional S.D. reduction in 907 Fig. 10. Additionally, a visualization of per-face texture S.D. 908 is provided in the supplementary material. 909

910 4.5.3 Hand Reconstruction with Extra In-the-wild Data

For hand reconstruction with extra in-the-wild data to fully exploit the advantage of our self-supervised method S²HAND(V), we use 34 train videos from YT 3D [40] and switch to orthogonal camera model with corresponding camera projection to enable learning with in-the-wild data without camera information.

⁹¹⁷ The results are presented in Table 8 and Fig. 11. Chang-⁹¹⁸ ing the camera model makes the model performance drop a

little, which may be caused by the ambiguity of camera focal 919 length. Then, imposing proposed quaternion loss boosts 920 the performance by 3.4%, which conforms to experiments 921 results in Section 4.5.2. Finally, adding extra in-the-wild data 922 further improves our model by 1.8%, resulting in 1.09cm in 923 MPVPE, which is the best result we can get on the HO-924 3D test set. From above, we see that our proposed method 925 is able to utilize extra in-the-wild data without any camera 926 information and benefits from these training data. It is worth 927 noticing that there is a certain domain gap between HO-928 3D and YT 3D since HO-3D regards hand-object interaction 929 scenario while YT 3D is mostly comprised of sign language 930 videos. However, the improvement on HO-3D has still been 931 witnessed, which we think proves the significance of utiliz-932 ing in the data and the advantage of the proposed method. 933

5 CONCLUSION

In this work, we have proposed self-supervised 3D hand 935 reconstruction models S²HAND and S²HAND(V) which 936 can be trained from a collection of hand images and video 937 sequences without manual annotations, respectively. The 938 3D hand reconstruction network in both models encodes the 939 input image into a set of meaningful semantic parameters 940 that represent hand pose, shape, texture, illumination, and 941 the camera viewpoint. These parameters can be decoded 942 into a textured 3D hand mesh as well as a set of 3D joints, 943 and in turn, the 3D mesh and joints can be projected into the 944 2D image space, which enables our network to be end-to-945 end learned. We further exploit the self-supervision signals 946

embedded in hand motion videos by developing a novel 947 quaternion loss and a texture and shape consistency loss 948 to obtain more accurate and consistent hand reconstruction. 949 Experimental results show that our models perform well 950 under noisy supervision sources captured from 2D hand 951 keypoint detection, and achieve comparable performance to 952 953 the state-of-the-art fully-supervised method. Moreover, the experiments on in-the-wild video data show that our self-954 supervised model is effective to learn useful information 955 from in-the-wild data to further improve its performance. 956

For the future study, we think the texture and shape 957 representation could be investigated deeply. Notice that we 958 only estimate a per-face color and adopt a low-resolution 959 hand model with 778 vertices. Besides, we rely on the shape 960 and texture regularization terms to learn hand appearance 961 from limited raw data under complex environments and po-962 tential hand-object interaction. These defects lead to the re-963 construction results having limited details. Recently, Corona 964 et al. [91] utilize implicit representation to learn hand shape 965 and texture of high quality. This points in a good direction 966 for further exploration. 967

ACKNOWLEDGMENT 968

This work was supported by the National Natural Sci-969 ence Foundation of China under Grant 62106177 and the 970 National Science Fund for Distinguished Young Scholars 971 of China under Grant 41725005. It was also supported 972 by the Joint Fund of the Ministry of Education of China 973 under Grant 8091B032156. The numerical calculation was 974 supported by the supercomputing system in the Super-975 computing Center of Wuhan University. 976

REFERENCES 977

- H. Lee, M. Billinghurst, and W. Woo, "Two-handed tangible in-978 [1] teraction techniques for composing augmented blocks," Virtual 979 Reality, vol. 15, no. 2, pp. 133-146, 2011. 980
- N. C. Camgoz, S. Hadfield, O. Koller, H. Ney, and R. Bowden, [2] 981 "Neural sign language translation," in Conference on Computer 982 983 Vision and Pattern Recognition, 2018.
- N. C. Camgoz, O. Koller, S. Hadfield, and R. Bowden, "Sign [3] 984 language transformers: Joint end-to-end sign language recognition 985 and translation," in Conference on Computer Vision and Pattern 986 987 Recognition, 2020.
- G. Garcia-Hernando, S. Yuan, S. Baek, and T.-K. Kim, "First-988 [4] person hand action benchmark with rgb-d videos and 3d hand 989 pose annotations," in Conference on Computer Vision and Pattern 990 Recognition, 2018. 991
- Z. Tu, J. Zhang, H. Li, Y. Chen, and J. Yuan, "Joint-bone fusion [5] 992 graph convolutional network for semi-supervised skeleton action 993 recognition," IEEE Transactions on Multimedia, 2022. 994
- M. Höll, M. Oberweger, C. Arth, and V. Lepetit, "Efficient physics-995 [6] based implementation for realistic hand-object interaction in vir-996 997 tual reality," in Conference on Virtual Reality and 3D User Interfaces, 2018. 998
- M. Parelli, K. Papadimitriou, G. Potamianos, G. Pavlakos, and [7] 999 P. Maragos, "Exploiting 3d hand pose estimation in deep learning-1000 based sign language recognition from rgb videos," in European 1001 1002 Conference on Computer Vision, 2020.
- Z. Tu, H. Li, D. Zhang, J. Dauwels, B. Li, and J. Yuan, "Action-stage [8] 1003 1004 emphasized spatiotemporal vlad for video action recognition, 1005 IEEE Transactions on Image Processing, vol. 28, no. 6, pp. 2799-2812, 2019. 1006
- L. Ge, H. Liang, J. Yuan, and D. Thalmann, "Robust 3d hand pose 1007 [9] estimation from single depth images using multi-view cnns," IEEE 1008 Transactions on Image Processing, vol. 27, no. 9, pp. 4422-4436, 1009 2018 1010

- [10] Z. Yu, J. S. Yoon, I. K. Lee, P. Venkatesh, J. Park, J. Yu, and 1011 H. S. Park, "Humbi: A large multiview dataset of human body 1012 expressions," in Conference on Computer Vision and Pattern 1013 Recognition, 2020. 1014
- [11] Z. Zhao, T. Wang, S. Xia, and Y. Wang, "Hand-3d-studio: 1015 A new multi-view system for 3d hand reconstruction," in IEEE International Conference on Acoustics, Speech, and Signal 1017 Processing, 2020.
- [12] G. Poier, D. Schinagl, and H. Bischof, "Learning pose specific representations by predicting different views," in Conference on Computer Vision and Pattern Recognition, 2018.
- [13] A. Armagan, G. Garcia-Hernando, S. Baek, S. Hampali, M. Rad, Z. Zhang, S. Xie, M. Chen, B. Zhang, F. Xiong et al., "Measuring generalisation to unseen viewpoints, articulations, shapes and objects for 3d hand pose estimation under hand-object interaction," in European Conference on Computer Vision, 2020.
- [14] Y. Chen, Z. Tu, L. Ge, D. Zhang, R. Chen, and J. Yuan, "So-handnet: Self-organizing network for 3d hand pose estimation with semisupervised learning," in International Conference on Computer Vision, 2019.
- [15] L. Ge, H. Liang, J. Yuan, and D. Thalmann, "Robust 3d hand pose estimation in single depth images: from single-view cnn to multi-view cnns," in Conference on Computer Vision and Pattern Recognition, 2016.
- [16] L. Huang, J. Tan, J. Liu, and J. Yuan, "Hand-transformer: Nonautoregressive structured modeling for 3d hand pose estimation," in European Conference on Computer Vision, 2020.
- [17] S. Yuan, G. Garcia-Hernando, B. Stenger, G. Moon, J. Yong Chang, K. Mu Lee, P. Molchanov, J. Kautz, S. Honari, L. Ge et al., "Depth-1039 based 3d hand pose estimation: From current achievements to future goals," in Conference on Computer Vision and Pattern 1041 Recognition, 2018.
- [18] V. Athitsos and S. Sclaroff, "Estimating 3d hand pose from a 1043 cluttered image," in Conference on Computer Vision and Pattern 1044 Recognition, 2003. 1045
- [19] Y. Cai, L. Ge, J. Cai, and J. Yuan, "Weakly-supervised 3d hand pose estimation from monocular rgb images," in European Conference on Computer Vision, 2018.
- [20] U. Iqbal, P. Molchanov, T. Breuel Juergen Gall, and J. Kautz, "Hand pose estimation via latent 2.5 d heatmap regression," in European 1050 Conference on Computer Vision, 2018.
- [21] A. Spurr, U. Iqbal, P. Molchanov, O. Hilliges, and J. Kautz, 1052 "Weakly supervised 3d hand pose estimation via biomechanical 1053 constraints," in European Conference on Computer Vision, 2020. 1054
- [22] C. Zimmermann and T. Brox, "Learning to estimate 3d hand pose 1055 from single rgb images," in International Conference on Computer 1056 Vision, 2017. 1057
- [23] L. Ge, Z. Ren, Y. Li, Z. Xue, Y. Wang, J. Cai, and J. Yuan, "3d hand 1058 shape and pose estimation from a single rgb image," in Conference 1059 on Computer Vision and Pattern Recognition, 2019. 1060
- [24] D. Kulon, R. A. Guler, I. Kokkinos, M. M. Bronstein, and 1061 S. Zafeiriou, "Weakly-supervised mesh-convolutional hand reconstruction in the wild," in Conference on Computer Vision and 1063 Pattern Recognition, 2020. 1065
- [25] Y. Hasson, B. Tekin, F. Bogo, I. Laptev, M. Pollefeys, and C. Schmid, "Leveraging photometric consistency over time for sparsely supervised hand-object reconstruction," in Conference on Computer Vision and Pattern Recognition, 2020.
- [26] Y. Hasson, G. Varol, D. Tzionas, I. Kalevatykh, M. J. Black, 1069 I. Laptev, and C. Schmid, "Learning joint reconstruction of hands and manipulated objects," in Conference on Computer Vision and 1070 1071 Pattern Recognition, 2019. 1072
- [27] N. Qian, J. Wang, F. Mueller, F. Bernard, V. Golyanik, and C. Theobalt, "Parametric hand texture model for 3d hand reconstruction and personalization," in European Conference on Computer Vision, 2020.
- [28] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Openpose: Realtime 1077 multi-person 2d pose estimation using part affinity fields." IEEE 1078 Transactions on Pattern Analysis and Machine Intelligence, vol. 43, 1079
- no. 01, pp. 172–186, 2019. Y. Chen, F. Wu, Z. Wang, Y. Song, Y. Ling, and L. Bao, "Self-[29] Transactions on Image Processing, vol. 29, pp. 8696-8705, 2020.
- [30] A. Tewari, M. Zollhofer, H. Kim, P. Garrido, F. Bernard, P. Perez, and C. Theobalt, "Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction," in International Conference on Computer Vision Workshops, 2017.

1016

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1040

1042

1046

1047

1048

1049

1051

1062

1064

1066

1067

1068

1073

1074

1075

1076

1080

1081

1082

1083

1084

1085

1086

- [31] Y. Chen, Z. Tu, D. Kang, L. Bao, Y. Zhang, X. Zhe, R. Chen, 1088 1089 and J. Yuan, "Model-based 3d hand reconstruction via selfsupervised learning," in Conference on Computer Vision and 1090 Pattern Recognition, 2021. 1091
- [32] Y. Cai, L. Ge, J. Cai, N. M. Thalmann, and J. Yuan, "3d hand pose 1092 estimation using synthetic data and weakly labeled rgb images, 1093 IEEE transactions on pattern analysis and machine intelligence, 1094 1095 vol. 43, no. 11, pp. 3739-3753, 2020.
- [33] L. Yang and A. Yao, "Disentangling latent hands for image synthe-1096 sis and pose estimation," in Conference on Computer Vision and 1097 Pattern Recognition, 2019. 1098
- K. Lin, L. Wang, and Z. Liu, "End-to-end human pose and mesh [34] 1099 reconstruction with transformers," in Conference on Computer 1100 Vision and Pattern Recognition, 2021. 1101
- P. Chen, Y. Chen, D. Yang, F. Wu, Q. Li, Q. Xia, and Y. Tan, "I2uv-[35] 1102 handnet: Image-to-uv prediction network for accurate and high-1103 fidelity 3d hand mesh modeling," in International Conference on 1104 Computer Vision, 2021. 1105
- S. Liu, H. Jiang, J. Xu, S. Liu, and X. Wang, "Semi-supervised 1106 [36] 3d hand-object poses estimation with interactions in time," in 1107 Conference on Computer Vision and Pattern Recognition, 2021. 1108
- Y. Zhou, M. Habermann, I. Habibie, A. Tewari, C. Theobalt, [37] 1109 and F. Xu, "Monocular real-time full body capture with inter-1110 part correlations," in Conference on Computer Vision and Pattern 1111 Recognition, 2021. 1112
- [38] X. Chen, Y. Liu, C. Ma, J. Chang, H. Wang, T. Chen, X. Guo, P. Wan, 1113 and W. Zheng, "Camera-space hand mesh recovery via semantic 1114 aggregation and adaptive 2d-1d registration," in Conference on 1115 Computer Vision and Pattern Recognition, 2021. 1116
- 1117 [39] Z. Cao, I. Radosavovic, A. Kanazawa, and J. Malik, "Reconstructing hand-object interactions in the wild," in International 1118 1119 Conference on Computer Vision, 2021.
- [40] D. Kulon, R. A. Guler, I. Kokkinos, M. M. Bronstein, and 1120 S. Zafeiriou, "Weakly-supervised mesh-convolutional hand recon-1121 struction in the wild," in Conference on Computer Vision and 1122 Pattern Recognition, 2020. 1123
- [41] I. Lim, A. Dielen, M. Campen, and L. Kobbelt, "A simple approach 1124 to intrinsic correspondence learning on unstructured 3d meshes,' 1125 in European Conference on Computer Vision, 2018. 1126
- [42] D. Kulon, H. Wang, R. A. Güler, M. M. Bronstein, and S. Zafeiriou, 1127 "Single image 3d hand reconstruction with mesh convolutions," 1128 in British Machine Vision Conference, 2019. 1129
- 1130 [43] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering,' 1131 Advances in neural information processing systems, vol. 29, p. 1132 1133 3844–3852, 2016.
- [44] L. Ge, Z. Ren, Y. Li, Z. Xue, Y. Wang, J. Cai, and J. Yuan, "3d hand 1134 shape and pose estimation from a single rgb image," in Conference 1135 on Computer Vision and Pattern Recognition, 2019. 1136
- J. Romero, D. Tzionas, and M. J. Black, "Embodied hands: Model-[45] 1137 ing and capturing hands and bodies together," ACM Transactions 1138 on Graphics (Proceedings of SIGGRAPH Asia), vol. 36, no. 6, 2017. 1139
- [46] A. Boukhayma, R. d. Bem, and P. H. Torr, "3d hand shape and 1140 pose from images in the wild," in Conference on Computer Vision 1141 and Pattern Recognition, 2019. 1142
- Y. Chen, Z. Tu, D. Kang, R. Chen, L. Bao, Z. Zhang, and 1143 [47]J. Yuan, "Joint hand-object 3d reconstruction from a single image 1144 with cross-branch feature fusion," IEEE Transactions on Image 1145 Processing, vol. 30, pp. 4008-4021, 2021. 1146
- C. Zimmermann, D. Ceylan, J. Yang, B. Russell, M. Argus, and 1147 [48] T. Brox, "Freihand: A dataset for markerless capture of hand pose 1148 and shape from single rgb images," in International Conference 1149 on Computer Vision, 2019. 1150
- S. Baek, K. I. Kim, and T.-K. Kim, "Pushing the envelope for rgb-[49] 1151 based dense 3d hand pose estimation via neural rendering, 1152 in Conference on Computer Vision and Pattern Recognition, 2019. 1153
- [50] "Weakly-supervised domain adaptation via gan and mesh 1154 model for estimating 3d hand poses interacting objects, 1155 Conference on Computer Vision and Pattern Recognition, 2020. 1156
- J. Wang, F. Mueller, F. Bernard, S. Sorli, O. Sotnychenko, N. Qian, [51] 1157 M. A. Otaduy, D. Casas, and C. Theobalt, "Rgb2hands: real-time 1158 1159 tracking of 3d hand interactions from monocular rgb video," ACM Transactions on Graphics, vol. 39, no. 6, pp. 1-16, 2020. 1160
- [52] 1161 Y. Zhou, M. Habermann, W. Xu, I. Habibie, C. Theobalt, and F. Xu, "Monocular real-time hand shape and motion capture using 1162 multi-modal data," in Conference on Computer Vision and Pattern 1163 Recognition, 2020. 1164

- [53] P. Panteleris, I. Oikonomidis, and A. Argyros, "Using a single rgb 1165 frame for real time 3d hand pose estimation in the wild," in Winter 1166 Conference on Applications of Computer Vision, 2018. 1167
- [54] S. Wu, C. Rupprecht, and A. Vedaldi, "Unsupervised learn-1168 ing of probably symmetric deformable 3d objects from images 1169 in the wild," in Conference on Computer Vision and Pattern 1170 Recognition, 2020. 1171 1172
- [55] S. Goel, A. Kanazawa, and J. Malik, "Shape and viewpoint without keypoints," in European Conference on Computer Vision, 2020.
- A. Spurr, A. Dahiya, X. Wang, X. Zhang, and O. Hilliges, "Self-[56] 1174 supervised 3d hand pose estimation from monocular rgb via 1175 contrastive learning," in International Conference on Computer 1176 Vision, 2021.
- S. Guo, E. Rigall, L. Qi, X. Dong, H. Li, and J. Dong, "Graph-based [57] 1178 cnns with self-supervised module for 3d hand pose estimation 1179 from monocular rgb," IEEE Transactions on Circuits and Systems 1180 for Video Technology, vol. 31, no. 4, pp. 1514-1525, 2020. 1181
- C. Wan, T. Probst, L. V. Gool, and A. Yao, "Self-supervised 3d [58] 1182 hand pose estimation through training by fitting," in Conference 1183 on Computer Vision and Pattern Recognition, 2019. 1184
- [59] A. Tewari, M. Zollhöfer, P. Garrido, F. Bernard, H. Kim, P. Pérez, and C. Theobalt, "Self-supervised multi-level face model learning for monocular reconstruction at over 250 hz," in Conference on Computer Vision and Pattern Recognition, 2018.
- [60] V. Blanz and T. Vetter, "A morphable model for the synthesis 1189 of 3d faces," in ACM Transactions on Graphics (Proceedings of 1190 SIGGRAPH), 1999, p. 187–194. 1191
- [61] M. de La Gorce, D. J. Fleet, and N. Paragios, "Model-based 3d 1192 hand pose estimation from monocular video," IEEE Transactions 1193 on Pattern Analysis and Machine Intelligence, vol. 33, no. 9, pp. 1194 1793-1805, 2011. 1195
- [62] M. de La Gorce, N. Paragios, and D. J. Fleet, "Model-based hand 1196 tracking with texture, shading and self-occlusions," in Conference 1197 on Computer Vision and Pattern Recognition, 2008.
- [63] Y. Cai, L. Ge, J. Liu, J. Cai, T.-J. Cham, J. Yuan, and N. M. Thalmann, "Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks," in International Conference on Computer Vision, 2019.
- [64] L. Chen, S.-Y. Lin, Y. Xie, Y.-Y. Lin, and X. Xie, "Temporal-aware self-supervised learning for 3d hand pose and mesh estimation in videos," in Winter Conference on Applications of Computer Vision, 2021.
- J. Yang, H. J. Chang, S. Lee, and N. Kwak, "Seqhand: Rgb-[65] sequence-based 3d hand pose and shape estimation," in European Conference on Computer Vision, 2020.
- [66] J. Zhang, Z. Tu, J. Yang, Y. Chen, and J. Yuan, "Mixste: Seq2seq mixed spatio-temporal encoder for 3d human pose estimation in video," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 13232-13242
- K. Fragkiadaki, S. Levine, P. Felsen, and J. Malik, "Recurrent net-[67] work models for human dynamics," in International Conference on Computer Vision, 2015.
- J. Martinez, M. J. Black, and J. Romero, "On human motion [68] prediction using recurrent neural networks," in Conference on Computer Vision and Pattern Recognition, 2017.
- [69] E. Aksan, M. Kaufmann, and O. Hilliges, "Structured prediction 1220 helps 3d human motion modelling," in International Conference 1221 on Computer Vision, 2019.
- A. Hernandez, J. Gall, and F. Moreno-Noguer, "Human mo-[70] 1223 tion prediction via spatio-temporal inpainting," in International 1224 Conference on Computer Vision, 2019. 1225
- [71] M. Tan and Q. V. Le, "Efficientnet: Rethinking model scaling for 1226 convolutional neural networks," in International Conference on 1227 Machine Learning, 2019. 1228
- H. Kato, Y. Ushiku, and T. Harada, "Neural 3d mesh renderer," in [72] Conference on Computer Vision and Pattern Recognition, 2018.
- A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in European Conference on Computer [73] Vision, 2016.
- X. Sun, B. Xiao, F. Wei, S. Liang, and Y. Wei, "Integral human pose [74] regression," in European Conference on Computer Vision, 2018.
- P. J. Huber, "Robust estimation of a location parameter," in Breakthroughs in statistics, 1992.
- Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity IEEE Transactions on Image Processing, vol. 13, no. 4, pp. 600-1240 612, 2004

1173

1185

1186

1187

1188

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1222

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

- [77] D. Pavllo, C. Feichtenhofer, M. Auli, and D. Grangier, "Mod-1242 1243 eling human motion with quaternion-based neural networks, International Journal of Computer Vision, vol. 128, no. 4, pp. 855-1244 872, 2020 1245
- [78] X. Zhang, S. Qin, Y. Xu, and H. Xu, "Quaternion product units for 1246 deep learning on 3d rotation groups," in Conference on Computer 1247 Vision and Pattern Recognition, 2020. 1248
- [79] K. Shoemake, "Animating rotation with quaternion curves," in 1249 ACM Transactions on Graphics (Proceedings of SIGGRAPH), 1250 vol. 19, no. 3, 1985, p. 245–254. 1251
- [80] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li, "On the continuity 1252 of rotation representations in neural networks," in Conference on 1253 Computer Vision and Pattern Recognition, 2019. 1254
- [81] S. Hampali, M. Rad, M. Oberweger, and V. Lepetit, "Honnotate: A 1255 method for 3d annotation of hand and object poses," in Conference 1256 on Computer Vision and Pattern Recognition, 2020. 1257
- J. Zhang, J. Jiao, M. Chen, L. Qu, X. Xu, and Q. Yang, "A hand [82] 1258 pose tracking benchmark from stereo matching," in International 1259 Conference on Image Processing, 2017. 1260
- [83] A. Knapitsch, J. Park, Q.-Y. Zhou, and V. Koltun, "Tanks and 1261 temples: Benchmarking large-scale scene reconstruction," ACM 1262 Transactions on Graphics, vol. 36, no. 4, pp. 1-13, 2017. 1263
- A. Kanazawa, J. Y. Zhang, P. Felsen, and J. Malik, "Learning 3d 1264 [84] human dynamics from video," in Conference on Computer Vision 1265 1266 and Pattern Recognition, 2019.
- A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, 1267 [85] Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differ-1268 entiation in pytorch," in Neural Information Processing Systems 1269 Workshops, 2017. 1270
- [86] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimiza-1271 tion," in International Conference for Learning Representations, 1272 1273 2014
- G. Moon and K. M. Lee, "I21-meshnet: Image-to-lixel prediction 1274 [87] network for accurate 3d human pose and mesh estimation from a 1275 single rgb image," in European Conference on Computer Vision, 1276 2020 1277
- [88] F. Mueller, F. Bernard, O. Sotnychenko, D. Mehta, S. Sridhar, 1278 D. Casas, and C. Theobalt, "Ganerated hands for real-time 3d hand 1279 tracking from monocular rgb," in Conference on Computer Vision 1280 and Pattern Recognition, 2018. 1281
- A. Spurr, J. Song, S. Park, and O. Hilliges, "Cross-modal deep [89] 1282 variational hand pose estimation," in Conference on Computer 1283 Vision and Pattern Recognition, 2018. 1284
- T. Theodoridis, T. Chatzis, V. Solachidis, K. Dimitropoulos, and [90] 1285 P. Daras, "Cross-modal variational alignment of latent spaces," 1286 in Conference on Computer Vision and Pattern Recognition 1287 Workshops, 2020. 1288
- E. Corona, T. Hodan, M. Vo, F. Moreno-Noguer, C. Sweeney, [91] 1289 R. Newcombe, and L. Ma, "Lisa: Learning implicit shape and appearance of hands," in <u>Conference on Computer Vision and</u> 1290 1291 Pattern Recognition, 2022. 1292



Zhigang Tu started his Master Degree in image processing at the School of Electronic Information, Wuhan University, China, 2008. In 2015, he received the Ph.D. degree in Computer Science from Utrecht University, Netherlands. From 2015 to 2016, he was a postdoctoral researcher at Arizona State University, US. Then from 2016 to 2018, he was a research fellow at the School of EEE, Nanyang Technological University, Singapore. He is currently a professor at the State Key Laboratory of Information Engineering in Sur-

veying, Mapping and Remote sensing, Wuhan University. His research interests include computer vision, image processing, video analytics, and machine learning. Special for motion estimation, video superresolution, object segmentation, action recognition and localization, and 1307 anomaly detection. He has co-/authored more than 50 articles on in-1308 ternational SCI-indexed journals and conferences. He is an Associate 1309 Editor of the SCI-indexed journal TVC (IF=1.456), a Guest Editor of 1310 JVCIR (IF=2.479) and Combinatorial Chemistry and High Throughput 1311 Screening (IF=1.195). He is the first organizer of the ACCV2020 Work-1312 shop on MMHAU (Japan). He received the "Best Student Paper" Award 1313 in the 4th Asian Conference on Artificial Intelligence Technology. 1314





Yujin Chen received the B.Eng. and M.Sc. in Geo-Information from Wuhan University. He is currently a Ph.D. student with the Visual Computing Lab, Technical University of Munich, Germany. His research interests include computer 1329 vision and machine learning. His current focus 1330 is on 3D scene understanding, pose estimation, 1331 motion analysis, and representation learning. 1332

Di Kang received the B.Eng. degree in Elec-1334 tronic Engineering and Information Science 1333 (EEIS) from the University of Science and Tech-1336 nology of China (USTC) in 2014, and the Ph.D. 1337 degree in computer science from City Univer-1338 sity of Hong Kong in 2019. In 2019, he jointed 1339 Tencent AI Lab, and is now serving as a senior 1340 researcher. His research interests include com-1341 puter vision, and deep learning. 1342

Linchao Bao received Ph.D. degree in Computer Science from City University of Hong Kong in 2015. Prior to that, he received M.S. degree in Pattern Recognition and Intelligent Systems from Huazhong University of Science and Technology in Wuhan, China. He was a research intern at Adobe Research from November 2013 to August 2014 and worked for DJI as an algorithm engineer from January 2015 to June 2016. His research interests include computer vision and graphics.

Bisheng Yang received the B.S. degree in en-1355 gineering survey, the M.S. and Ph.D. degrees 1356 in photogrammetry and remote sensing from 1357 Wuhan University, Wuhan, China, in 1996, 1999, 1358 and 2002, respectively. From 2002 to 2006, 1359 he was a Post-doctoral Research Fellow with 1360 the University of Zurich, Zurich, Switzerland. 1361 Since 2007, he has been a Professor with the 1362 State Key Laboratory of Information Engineering 1363 in Surveying, Mapping, and Remote Sensing, 1364 Wuhan University, where he is currently the Di-1365

rector of the 3S and Network Communication Laboratory. He has hosted a project of the National High Technology Research and Development 1367 Program, a key project of the Ministry of Education, and four National Scientific Research Foundation Projects of China. His current research 1369 interests include 3-D geographic information systems, urban modeling, 1370 and digital city. He was a Guest Editor of the ISPRS Journal of Photogrammetry and Remote Sensing and Computers & Geosciences. Junsong Yuan is Professor and Director of Vi-



1373 sual Computing Lab at Department of Computer 1374 Science and Engineering (CSE), State Univer-1375 sity of New York at Buffalo, USA. Before join-1376 ing SUNY Buffalo, he was Associate Professor 1377 (2015-2018) and Nanyang Assistant Professor 1378 (2009-2015) at Nanyang Technological Univer-1379 sity (NTU), Singapore. He obtained his Ph.D. 1380 from Northwestern University in 2009, M.Eng. 1381 from National University of Singapore in 2005, 1382 and B.Eng. from Huazhong University of Science 1383

Technology (HUST) in 2002. His research interests include computer 1384 vision, pattern recognition, video analytics, human action and gesture 1385 analysis, large-scale visual search and mining. He received Best Paper 1386 Award from IEEE Trans. on Multimedia, Nanyang Assistant Profes-1387 sorship from NTU, and Outstanding EECS Ph.D. Thesis award from 1388 Northwestern University. He served as Associate Editor of IEEE Trans. 1389 on Image Processing (T-IP), IEEE Trans. on Circuits and Systems for 1390 Video Technology (T-CSVT), Machine Vision and Applications (MVA), 1391 and Senior Area Editor of Journal of Visual Communications and Image 1392 Representation (JVCI). He was Program Co-Chair of IEEE Conf. on 1393 Multimedia Expo (ICME'18), and served as Area Chair for CVPR, ICCV, 1394 and ACM MM. He is a Fellow of IEEE and IAPR. 1395

16

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1366

1368

1371

1372

1315

1316