

Supplementary Material for Extract-and-Adaptation Network for 3D Interacting Hand Mesh Recovery

JoonKyu Park^{1*} Daniel Sungho Jung^{2,3*} Gyeongsik Moon^{4*} Kyoung Mu Lee^{1,2,3}

¹Dept. of ECE&ASRI, ²IPAI, Seoul National University, Korea

³SNU-LG AI Research Center, ⁴Meta Reality Labs Research

{jkipark0825, dqj5182}@snu.ac.kr, mks0601@meta.com, kyoungmu@snu.ac.kr

In this supplementary material, we provide more experiments, discussions, and other details that could not be included in the main text due to the lack of pages. The contents are summarized below:

- S1. Detailed architecture
 - S1.1. FuseFormer
 - S1.2. EABlock
- S2. Number of adaptation stages in EABlock
- S3. MPJPE comparison
- S4. Measuring the pose difference
- S5. Additional qualitative comparisons
- S6. Discussion
 - S6.1. Limitations
 - S6.2. Societal impacts

S1. Detailed architecture

In the main manuscript, we introduce the Extract-and-Adaptation Network (EANet) as a means of effectively extracting and processing interactions between two hands. The proposed EANet mainly operates upon the key component, EABlock, to extract and adapt features from the two hands. In this section, we provide a detailed explanation of the EABlock and its primary component, the FuseFormer, which are illustrated in Figures 4 and 5 of the main manuscript, respectively.

S1.1. FuseFormer

In Algorithm S1, we show the inference process of FuseFormer in the extract stage of EABlock. From two hand features (\mathbf{F}_L and \mathbf{F}_R), JoinToken \mathbf{t}_J is obtained by passing a concatenated two hand features to a fully connected layer after reshaping. To obtain SimToken \mathbf{t}_S , we first reshape each hand feature (\mathbf{F}_L or \mathbf{F}_R) and concatenate the reshaped left and right hand features with a class token \mathbf{t}_{cls} . A concatenated token \mathbf{t} is made after reshaping the concatenated feature of \mathbf{t}_{cls} , \mathbf{F}_L , and \mathbf{F}_R . The concatenated token \mathbf{t} is then used for extracting query \mathbf{q}_{SA} , key \mathbf{k}_{SA} , and value \mathbf{v}_{SA} with separate linear layers in a self-attention (SA) based Trans-

former. After preparing two tokens (\mathbf{t}_J and \mathbf{t}_S), a cross-attention (CA) based Transformer processes the two novel tokens, JoinToken \mathbf{t}_J as query \mathbf{q}_{CA} and SimToken \mathbf{t}_S as key-value pair (\mathbf{k}_{CA} and \mathbf{v}_{CA}), to effectively fuse two hand features (\mathbf{F}_L and \mathbf{F}_R) and output an interaction feature \mathbf{F}_{inter} .

S1.2. EABlock

We show the inference process of EABlock in Algorithm S2. From left hand feature \mathbf{F}_L and right hand feature \mathbf{F}_R , a FuseFormer in the extract stage of EABlock first extracts an interaction feature \mathbf{F}_{inter} . Then, with the extracted interaction feature \mathbf{F}_{inter} in the extract stage, each FuseFormer in the adaptation stage of EABlock fuses each hand feature (\mathbf{F}_L or \mathbf{F}_R) and the interaction feature \mathbf{F}_{inter} to produce an adapted interaction feature for each hand (\mathbf{F}_{interL} or \mathbf{F}_{interR}). Lastly, our EABlock constructs the final outputs (\mathbf{F}_L^* and \mathbf{F}_R^*) after a fully connected layer on each adapted interaction feature (\mathbf{F}_{interL} or \mathbf{F}_{interR}) and then a concatenation to each hand feature (\mathbf{F}_L or \mathbf{F}_R), respectively.

S2. Number of adaptation stages in EABlock

To adapt the extracted interaction feature to each hand feature, our EABlock fuses an interaction feature and one hand feature (*i.e.* left hand feature or right hand feature) to obtain an adapted interaction feature for each hand. Here, we justify our strategy to adapt only once to each hand feature. Table S1 shows that a single iteration of the adaptation stage shows comparable results to two iterations of the adaptation stages. Considering the additional costs followed by additional adaptation stages, we set the number of adaptation stages to 1.

S3. MPJPE comparison with state-of-the-art methods

Table 5 in the main manuscript provides a comparison of MPVPE and MRRPE for various methods, while Table S2 presents additional results based on MPJPE. We follow the

Algorithm S1 Pseudocode of FuseFormer in a PyTorch-style

```

1: class FuseFormer(nn.Module):
2:     def __init__():
3:         FC = nn.Linear(1024, 512)
4:         t_cls = nn.Parameter(512, 1)
5:         Q_SA = nn.Linear(512, 512)
6:         K_SA = nn.Linear(512, 512)
7:         V_SA = nn.Linear(512, 512)
8:         Q_CA = nn.Linear(512, 512)
9:         K_CA = nn.Linear(512, 512)
10:        V_CA = nn.Linear(512, 512)
11:         $\psi : \mathbb{R}^{1024 \times 8 \times 8} \rightarrow \mathbb{R}^{64 \times 1024}$  # reshape
12:         $\phi : \mathbb{R}^{512 \times 8 \times 8} \rightarrow \mathbb{R}^{512 \times 64}$  # reshape
13:         $\eta : \mathbb{R}^{512 \times 64} \rightarrow \mathbb{R}^{129 \times 512}$  # reshape
14:        softmax = nn.Softmax(dim=-1)
15:        MLP = nn.Sequential(
16:            nn.Linear(512, 512*4),
17:            nn.Linear(512*4, 512))
18:        MLP' = nn.Sequential(
19:            nn.Linear(512, 512*4),
20:            nn.Linear(512*4, 512))
21:         $d_{k_{SA}}, d_{k_{CA}} = 512, 512$ 
22:        def forward( $\mathbf{F}_L, \mathbf{F}_R$ ): #  $\mathbf{F}_L$  and  $\mathbf{F}_R \in \mathbb{R}^{512 \times 8 \times 8}$ 
23:            # get JoinToken
24:             $\mathbf{t}_J = \text{FC}(\psi(\text{torch.cat}((\mathbf{F}_L, \mathbf{F}_R)))) \in \mathbb{R}^{64 \times 512}$ 
25:            # get SimToken
26:             $\mathbf{F}_L = \phi(\mathbf{F}_L) \in \mathbb{R}^{512 \times 64}$ 
27:             $\mathbf{F}_R = \phi(\mathbf{F}_R) \in \mathbb{R}^{512 \times 64}$ 
28:             $\mathbf{t} = \eta(\text{torch.cat}((\mathbf{t}_{cls}, \mathbf{F}_L, \mathbf{F}_R))) \in \mathbb{R}^{129 \times 512}$ 
29:             $\mathbf{q}_{SA}, \mathbf{k}_{SA}, \mathbf{v}_{SA} = \text{Q\_SA}(\mathbf{t}), \text{K\_SA}(\mathbf{t}), \text{V\_SA}(\mathbf{t})$ 
30:             $\mathbf{r} = \text{softmax}((\mathbf{q}_{SA} @ \mathbf{k}_{SA}^T) / \sqrt{d_{k_{SA}}}) \mathbf{v}_{SA} + \mathbf{t}$ 
31:             $\mathbf{t}_S = \mathbf{r} + \text{MLP}(\mathbf{r})$ 
32:            # process two tokens
33:             $\mathbf{q}_{CA} = \text{Q\_CA}(\mathbf{t}_J)$ 
34:             $\mathbf{k}_{CA}, \mathbf{v}_{CA} = \text{K\_CA}(\mathbf{t}_S), \text{V\_CA}(\mathbf{t}_S)$ 
35:             $\mathbf{r}' = \text{softmax}((\mathbf{q}_{CA} @ \mathbf{k}_{CA}^T) / \sqrt{d_{k_{CA}}}) \mathbf{v}_{CA} + \mathbf{t}_J$ 
36:             $\mathbf{F}_{inter} = \mathbf{r}' + \text{MLP}'(\mathbf{r}')$ 
37:            return  $\mathbf{F}_{inter}$ 

```

# of adaptation	MPJPE			MPVPE			MRRPE
	Single	Two	All	Single	Two	All	
0	10.66	11.81	11.25	8.37	11.24	9.61	21.00
1	9.62	11.54	10.58	7.86	10.78	9.13	18.82
2	9.60	11.90	10.75	8.14	10.66	9.19	18.82

Table S1: Comparison of models with various number of adaptation stages.

evaluation protocol of prior studies by using a pre-defined joint regression matrix to obtain the joint positions from the estimated meshes, which are scaled using the ground truth bone length. Note that the MPJPE results for Keypoint

Algorithm S2 Pseudocode of EABlock in a PyTorch-style

```

1: class EABlock(nn.Module):
2:     def __init__():
3:         FuseFormer_extract = FuseFormer()
4:         FuseFormer_L = FuseFormer()
5:         FuseFormer_R = FuseFormer()
6:         FC = nn.Linear(512, 128)
7:     def forward( $\mathbf{F}_L, \mathbf{F}_R$ ): #  $\mathbf{F}_L$  and  $\mathbf{F}_R \in \mathbb{R}^{512 \times 8 \times 8}$ 
8:         # interaction feature extract
9:          $\mathbf{F}_{inter} = \text{FuseFormer\_extract}(\mathbf{F}_L, \mathbf{F}_R)$ 
10:        # interaction feature adaptation
11:         $\mathbf{F}_{interL} = \text{FuseFormer\_L}(\mathbf{F}_L, \mathbf{F}_{inter})$ 
12:         $\mathbf{F}_{interR} = \text{FuseFormer\_R}(\mathbf{F}_R, \mathbf{F}_{inter})$ 
13:         $\mathbf{F}_L^* = \text{torch.cat}((\mathbf{F}_L, \text{FC}(\mathbf{F}_{interL})))$ 
14:         $\mathbf{F}_R^* = \text{torch.cat}((\mathbf{F}_R, \text{FC}(\mathbf{F}_{interR})))$ 
15:        return  $\mathbf{F}_L^*, \mathbf{F}_R^*$ 

```

Transformer in Table S2 differ from those reported in their main manuscript, as our results are obtained using their official codes and pre-trained weights based on the mesh representation, while their reported numbers were based on a 2.5D representation. Our results demonstrate that EANet outperforms other methods in 3D joint estimation, indicating its effectiveness.

Methods	MPJPE		
	Single	Two	All
Zhang <i>et al.</i> [8]	-	13.48	-
Keypoint Transformer [2]	11.04	16.22	13.82
IntagHand [3]	-	8.79	-
EANet (Ours)	5.19	6.57	5.88

Table S2: Quantitative comparison with state-of-the-art methods on InterHand2.6M [4] dataset.

S4. Measuring the pose difference for Figure 9

In L833-L850 and Figure 9, we present the pose difference between two hands across different sequences of the InterHand2.6M dataset. To compute the pose difference for each selected sequence, we begin by flipping a right hand to a left hand. Next, we subtract the root joint (wrist) position of each hand. Finally, we align the global rotation of the two hands to solely determine the finger pose symmetry.

S5. Additional qualitative comparisons

In the below figures, we further show the qualitative comparisons with other state-of-the-art methods [3, 2]. Specifically, in Figure S1 and S2, we show visual comparisons on InterHand2.6M [4] dataset. In Figure S3, we show visual comparisons on in-the-wild interacting hand images.

S6. Discussion

S6.1. Limitations

Despite recent progress on domain adaptation [1, 5, 6] that provides several techniques for reducing the discrepancy between training data and test data, our work has not yet employed such designs for better generalization ability. In Figure S3, our EANet shows reasonable results on in-the-wild images without such designs. However, we believe that additional designs to improve generalization ability on in-the-wild images can further improve the performance of EANet in future works.

S6.2. Societal impacts

Our EANet performs robust 3D interacting hand mesh recovery with potential applications on technologies connecting real-world and virtual-world such as AR and VR. Especially, the work is useful in cases where two hands are often interacting, including virtual meetings and virtual events.

License of the Used Assets

- [InterHand2.6M dataset \[4\]](#) is CC-BY-NC 4.0 licensed.
- [HIC dataset \[7\]](#) is publicly available dataset.
- [Zhang *et al.* \[8\] codes](#) are released for academic research only, and it is free to researchers from educational or research institutes for non-commercial purposes.
- [Keypoint Transformer \[2\] codes](#) are released for academic research only, and it is free to researchers from educational or research institutes for non-commercial purposes.
- [IntagHand \[3\] codes](#) are released for academic research only, and it is free to researchers from educational or research institutes for non-commercial purposes.

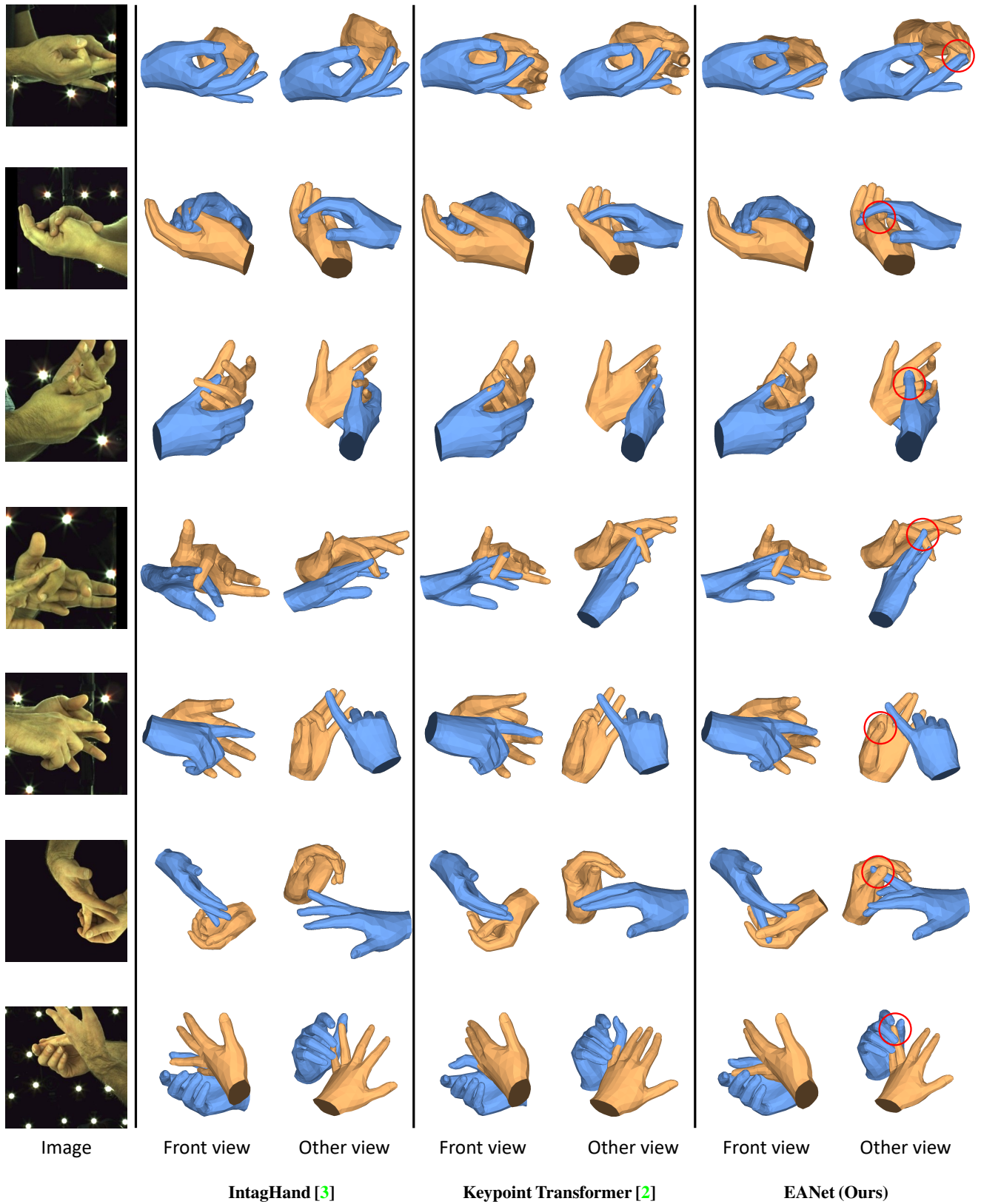


Figure S1: Visual comparison with state-of-the-art methods on InterHand2.6M [4]. The red circles highlight regions where our EANet is correct, while others are wrong.

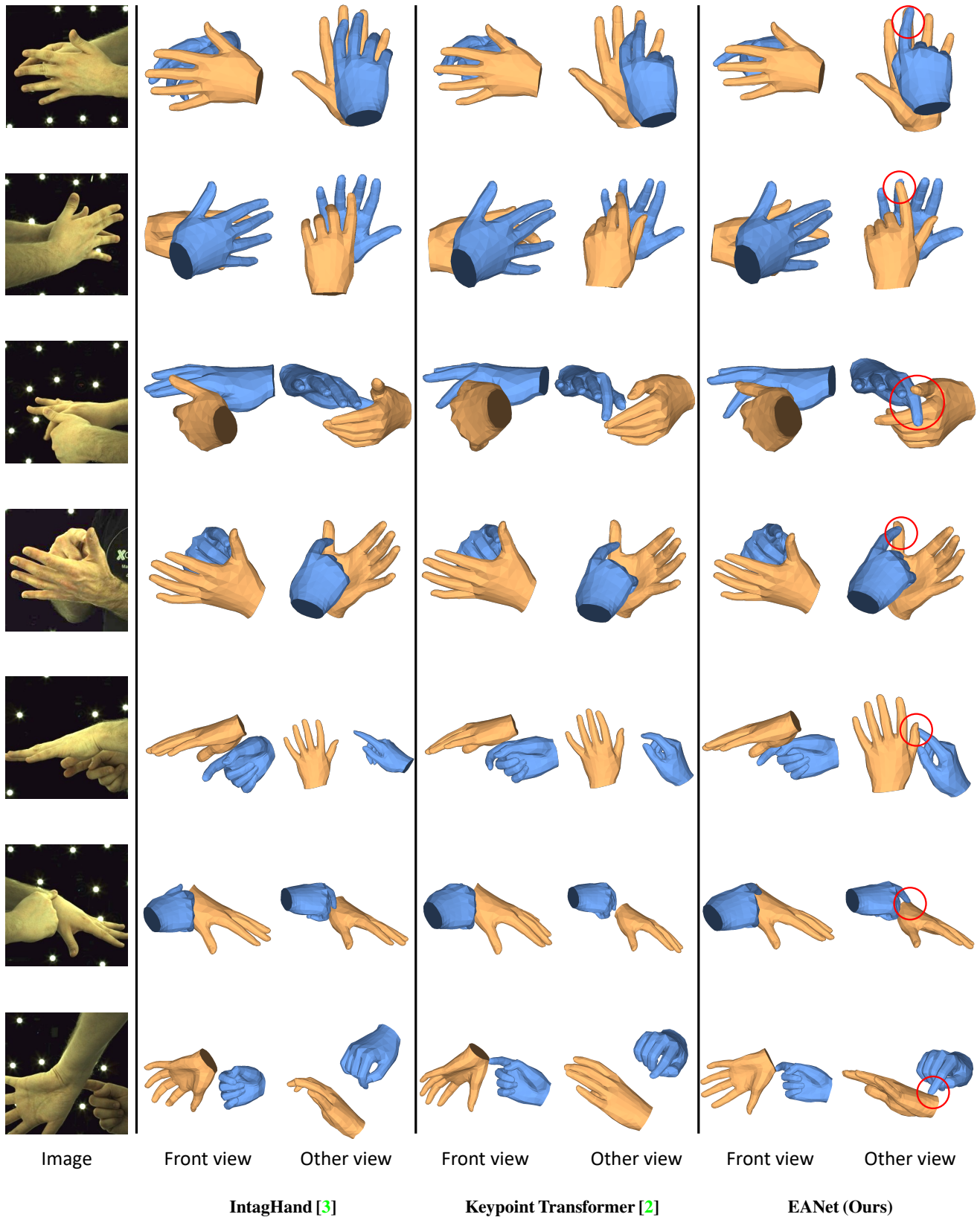


Figure S2: Visual comparison with state-of-the-art methods on InterHand2.6M [4]. The red circles highlight regions where our EANet is correct, while others are wrong.

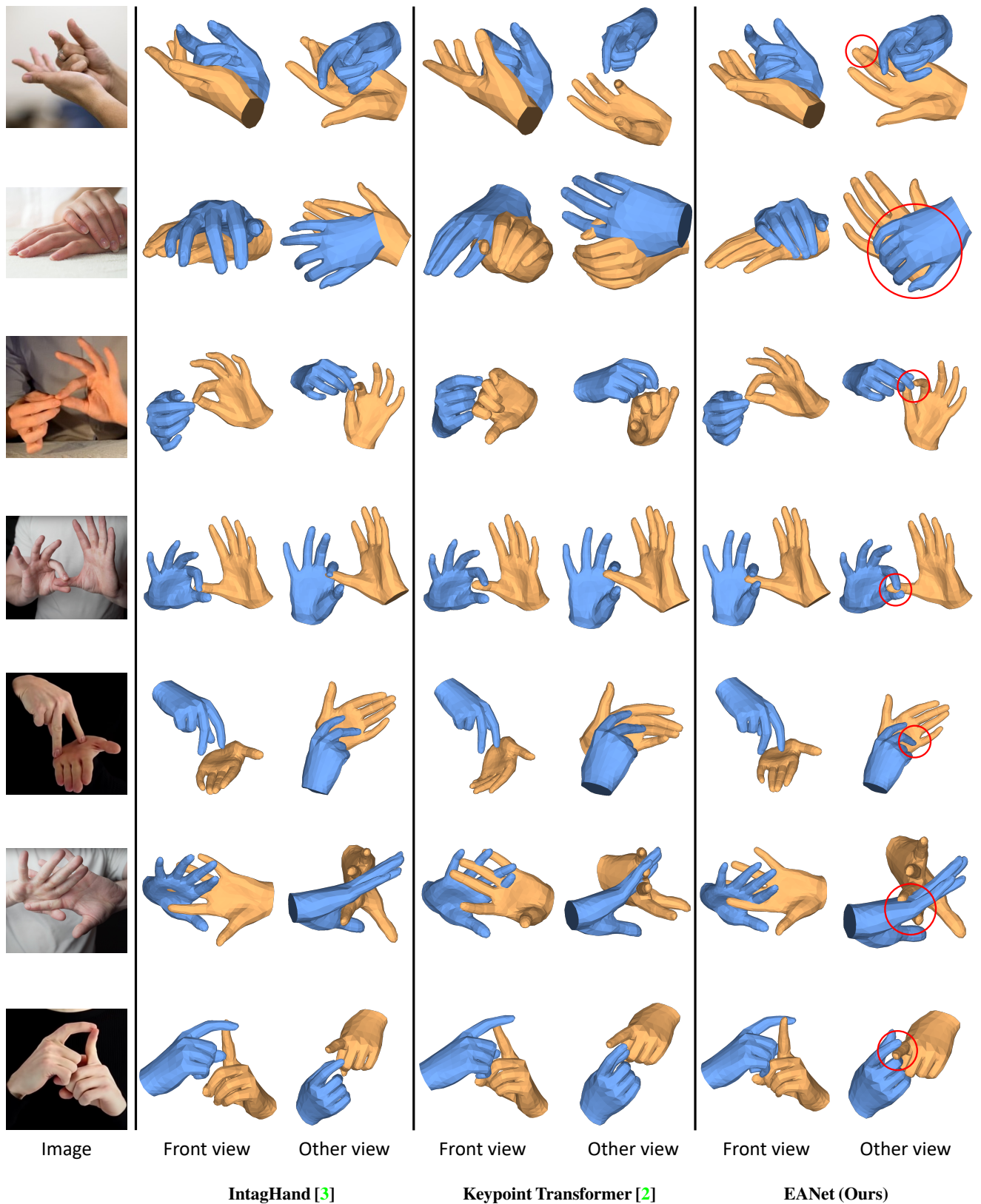


Figure S3: Visual comparison with state-of-the-art methods on in-the-wild images. The red circles highlight regions where our EANet is correct, while others are wrong. The images are crawled from Google and YouTube.

References

- [1] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, 2015. 3
- [2] Shreyas Hampali, Sayan Deb Sarkar, Mahdi Rad, and Vincent Lepetit. Keypoint Transformer: Solving joint identification in challenging hands and object interactions for accurate 3D pose estimation. In *CVPR*, 2022. 2, 3, 4, 5, 6
- [3] Mengcheng Li, Liang An, Hongwen Zhang, Lianpeng Wu, Feng Chen, Tao Yu, and Yebin Liu. Interacting attention graph for single image two-hand reconstruction. In *CVPR*, 2022. 2, 3, 4, 5, 6
- [4] Gyeongsik Moon, Shoou-I Yu, He Wen, Takaaki Shiratori, and Kyoung Mu Lee. InterHand2.6M: A dataset and baseline for 3D interacting hand pose estimation from a single RGB image. In *ECCV*, 2020. 2, 3, 4, 5
- [5] Zhongyi Pei, Zhangjie Cao, Mingsheng Long, and Jianmin Wang. Multi-adversarial domain adaptation. In *AAAI*, 2018. 3
- [6] Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation. In *ICLR*, 2017. 3
- [7] Dimitrios Tzionas, Luca Ballan, Abhilash Srikantha, Pablo Aponte, Marc Pollefeys, and Juergen Gall. Capturing hands in action using discriminative salient points and physics simulation. *IJCV*, 2016. 3
- [8] Baowen Zhang, Yangang Wang, Xiaoming Deng, Yinda Zhang, Ping Tan, Cuixia Ma, and Hongan Wang. Interacting two-hand 3D pose and shape reconstruction from single color image. In *ICCV*, 2021. 2, 3