# Keypoint Transformer: Solving Joint Identification in Challenging Hands and Object Interactions for Accurate 3D Pose Estimation

Shreyas Hampali[(1)], Sayan Deb Sarkar[(1)], Mahdi Rad[(1)], Vincent Lepetit[(2,1)]

[(1)]Institute for Computer Graphics and Vision, Graz University of Technology, Graz, Austria
[(2)]LIGM, Ecole des Ponts, Univ Gustave Eiffel, CNRS, Marne-la-Vallée, France

`{<firstname>.<lastname>}@icg.tugraz.at, vincent.lepetit@enpc.fr`
Project Website: `https://www.tugraz.at/index.php?id=57823`

## Abstract

*We propose a robust and accurate method for estimating the 3D poses of two hands in close interaction from a single color image. This is a very challenging problem, as large occlusions and many confusions between the joints may happen. State-of-the-art methods solve this problem by regressing a heatmap for each joint, which requires solving two problems simultaneously: localizing the joints and recognizing them. In this work, we propose to separate these tasks by relying on a CNN to first localize joints as 2D keypoints, and on self-attention between the CNN features at these keypoints to associate them with the corresponding hand joint. The resulting architecture, which we call "Keypoint Transformer", is highly efficient as it achieves state-of-the-art performance with roughly half the number of model parameters on the InterHand2.6M dataset. We also show it can be easily extended to estimate the 3D pose of an object manipulated by one or two hands with high performance. Moreover, we created a new dataset of more than 75,000 images of two hands manipulating an object fully annotated in 3D and will make it publicly available.*

## 1. Introduction

3D hand pose estimation has the potential to make virtual reality, augmented reality, and interaction with computers and robots much more intuitive. Recently, significant progress has been made for single-hand pose estimation using depth maps and even single RGB images. Being able to deal with RGB images is particularly attractive as it does not require a power-hungry active sensor. Many approaches have been proposed, mostly based on direct prediction with different convolutional network architectures [17, 21, 36, 43, 52, 57, 70] of the 3D joint locations or angles, or relying on rendering for fine pose estimation and tracking [2, 14, 39, 48, 58].
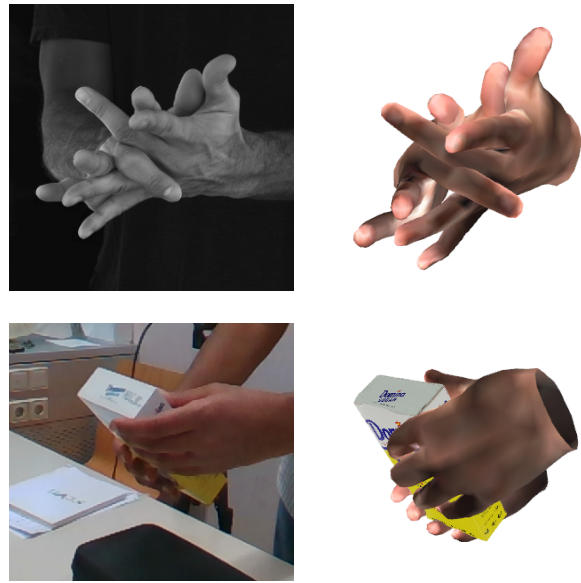


Figure 1. Our approach accurately predicts 3D hand and object poses from a single RGB image in challenging scenarios including complex hand interactions (top) and 2 hands interacting with an object where the hands can be severely occluded (bottom). The bottom example is from the $H_2O$-3D dataset we also introduce in this paper, which contains challenging, fully and accurately 3D-annotated, video sequences of two hands manipulating objects.

In contrast to single-hand pose estimation, two-hand pose estimation has received much less attention. This problem is indeed significantly harder: The appearance similarities between the joints of the two hands make their identification extremely challenging. Moreover, in close interaction, some of the joints of a hand are likely to be occluded by the other hand or itself. Thus, first detecting the left and right hands and then independently predicting their 3D poses [15, 43] performs poorly in close interaction scenarios. Bottom-up approaches [37, 62] directly estimate the 2D joint locations and their depths using one heatmap per joint. However, as shown in Fig. 2, the similarity in appearances
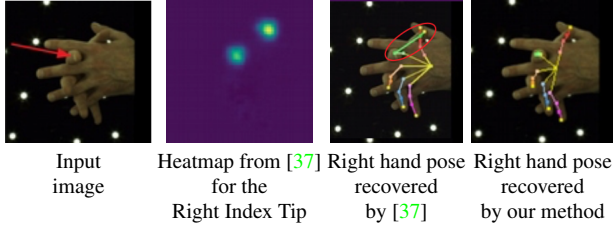
| Input image | Heatmap from [37] for the Right Index Tip | Right hand pose recovered by [37] | Right hand pose recovered by our method |

Figure 2. **Similar appearances between joints and partial occlusions make previous methods prone to failure.** The Inter-Net state-of-the-art method [37] predicts a heatmap for each joint but the predicted heatmaps can become ambiguous, resulting in failures when predicting the hand pose (the hand on the back in this example). Our approach explicitly models the relationship between keypoints resulting in more accurate poses. More examples can be found in the supplementary material.

of the joints and severe occlusions degrade the quality of heatmaps failing to localize the joints accurately. More recent works [12,25,66] have attempted to alleviate this problem by exploiting joint-segmentation, joint-visibility, or by adding more refinement layers increasing the overall complexity of the network. By exploiting only keypoints, our method outperforms these methods by a large margin with a significantly smaller model.

As shown in Fig. 3, instead of aiming to localize and recognize the hand joints simultaneously, we estimate the 3D poses of the hands in three stages: (1) We first detect "keypoints", which are potential joint locations in the image, by predicting a *single* heatmap. These keypoints do not have to exactly match all the hand joints: The 3D poses we predict are still correct if some joints are not detected as keypoints, and if some keypoints do not correspond to joints. (2) Then, we associate the keypoints with the corresponding joint or to the background in the case of false positives, on the basis of the keypoint locations and their image features. This is done for all the keypoints simultaneously to exploit mutual constraints, using the self-attention mechanism. (3) Finally, we predict the 3D hand poses using a cross-attention module, which selects keypoints associated with each of the hand joints. Our approach is agnostic to the parameterization of the pose and we consider three different hand pose representations.

Our architecture, which we call "Keypoint Transformer", is therefore designed to explicitly disambiguate the identity of the keypoints and performs very well even on complex configurations. Fig. 1 shows its output on two challenging examples, using the MANO [49] mesh as the output representation. Our architecture is related to the "Detection Transformer" (DETR) [8] architecture. DETR uses all the spatial features from a low-resolution CNN feature map, combined with learned location queries to detect objects in an image. The high computational complexity of the Transformer restricts DETR from using higher resolution CNN feature maps. As we show in our experiments,

using the DETR-style architecture for hand pose estimation results in lower accuracy and we hypothesize that this is due to the use of lower resolution feature maps and features from the entire image.

We train and evaluate our architecture on the recent InterHand2.6M hand-hand [37] and HO-3D hand-object [14] interaction datasets. We also introduce a challenging dataset of videos with two hands interacting with an object with complete and accurate 3D annotations without markers. This dataset is based on the work of [14], and we call it $H_2O$-3D. Our method achieves state-of-the-art performance on existing hand-interaction datasets and serves as a strong baseline for the $H_2O$-3D dataset. Our experiments show that on InterHand2.6M, our method achieves state-of-the-art performance with roughly half the number of model parameters. We carry out several ablation studies and compare with strong baselines to prove the efficacy of our approach.

## 2. Related Work

Many approaches have already been proposed for hand or object pose estimation from either RGB images or depth maps. Here we focus mainly on works that estimate hand poses during hand-hand or hand-object interactions. We also discuss recent advances in Transformer architectures for computer vision as they are highly relevant to our work.

### 2.1. Interacting Hand Pose Estimation

Hand pose estimation methods can be broadly classified as generative, discriminative, or hybrid approaches. Generative methods [14, 30, 40–42, 60] fit a parametric hand model to an observed image or depth map by minimizing a fitting error under some constraints. Discriminative methods [5, 16, 17, 23, 37, 43, 57, 71] mostly directly predict the hand pose from a single frame. Generative methods often rely heavily on tracking and are prone to drift whereas discriminative methods tend to generalize poorly to unseen images [1]. Hybrid approaches [4, 7, 15, 38, 51, 53, 55, 56, 59, 62, 63] try to combine the best of these two worlds by using discriminative methods to detect visual cues in the image followed by model fitting.

Earlier methods [30, 40, 41] for generative hand pose estimation during interaction used complex optimization methods to fit a parametric hand model to RGBD data. [14] proposed multi-frame optimization to fit hand and object models to RGBD data from multiple RGBD cameras. Generative methods alone often lose tracking during close interactions or occlusions and are hence combined with discriminative methods to guide the optimization.

[4, 59] detect the fingertips as discriminative points and used them in the optimization along with a collision term and physical modelling. Recently, [51] proposed high-fidelity hand surface tracking of hand-hand interactions in

a multi-view setup where the regressed 3D hand joint locations were used for initializing the tracking. [7,15,38,43,62] compute dense features or keypoints from a single RGB or depth image and fit a hand model [49] to these estimates with physical constraints and joint angle constraints. Fully discriminative methods [16, 17, 37, 57] jointly estimate the 3D joint locations or hand model parameters of both the interacting hands or the interacting hand and the object by incorporating contacts and inter-penetrations in the training. [23] estimates the hand-object surface using implicit representation that naturally allows modelling of the contact regions between hand and object. [12, 25] improve the accuracy of 3D pose estimation during hand-hand interaction scenarios by incorporating joint-visibility and part-segmentation cues, whereas [66] utilize refinement layers to iteratively refine the estimated poses.

By contrast with the above mentioned approaches designed specifically for hand-hand or hand-object interaction scenarios, we propose in this work a unified discriminative approach for all hand interaction scenarios. Further, many previous discriminative methods perform poorly during close hand interactions due to similarity in appearance of the joints. In this work, we model relationship between all detected joints in the image resulting in more accurate pose estimation while keeping the model complexity low.

The success of discriminative methods depend on the variability of training data and several hand interaction datasets have been proposed. [13] first provided a marker-based hand-object interaction dataset using RGBD cameras. [14, 71] and [17] respectively proposed real and synthetic hand-object interaction dataset with a single hand manipulating an object, while [29] recently developed a two-hands and object interaction dataset. [5] proposed single and two-hand object interaction dataset using infrared camera for contact annotations. [37] developed a large-scale two-hand interaction dataset using semi-automatic annotation process with many close interactions. [54] used MoCap to obtain pose of full body, hand, and object during interaction and used it to generate realistic grasp on unseen objects.

In this work, we also introduce a challenging two-hands-and-object interaction dataset which we created using the optimization method of [14]. Our dataset is made of videos of two hands from different subjects manipulating an object from the YCB dataset [64], annotated with the 3D poses of the hands and the object. Our architecture already performs well on this dataset and constitutes a strong baseline.

### 2.2. Transformers in Computer Vision

Transformers have recently been increasingly gaining popularity for vision related problems [24]. Features are often extracted from a CNN backbone and different architectures have been proposed to solve object detection [8, 69], image classification [11], pose estimation [6, 20, 32, 33] and

low-level image tasks [28, 65]. We refer the reader to [24] for a detailed survey.

[8, 69] proposed to combine a CNN backbone with a Transformer to detect objects in an image. [32] proposed to reconstruct the vertices of a single human body or hand from an RGB image using multiple Transformer encoder layers and achieved state-of-the-art performance. [33] improved [32] by using graph convolutions along with a Transformer encoder. [20] estimated a 3D pose from hand point-cloud data using a Transformer encoder-decoder architecture and proposed to generate query embeddings from input point-cloud instead of learning them as in [8, 69]. While these works are aimed at single hand pose estimation and their extension to two hands is non-trivial, our architecture is designed to estimate single and two-hands poses along with the object pose during hand-object interaction from the input RGB image.

In a closely related work, [6] solves the multi-person 2D pose estimation problem using detected keypoints and person centers, by associating the joint keypoints to the correct person center using attention. There are however several key differences: In our case, the hand centers get very close to each other during close interaction, and the approach in [6] would not be transferable. More importantly, hand joints are much more ambiguous than "body joints" as they look very similar to each other. Our method is also robust to undetected and falsely detected keypoints as we show in our discussions, while [6] cannot handle undetected keypoints. Further, we show that, by randomly sampling keypoints on the object, we can easily extend our method to 3D object pose estimation during hand-object interactions.

## 3. Method

As shown in Fig. 3, our architecture first detects keypoints that are likely to correspond to the 2D locations of hand joints and encodes them as input to the keypoint-joint association stage. The keypoints are encoded with their spatial locations and the image features at these locations. The self-attention layers in the Keypoint Transformer disambiguate the keypoints and associates them with different joint types and a background class. The (single) cross-attention layer then selects these "identity-aware keypoints" to predict root-joint-relative pose parameters of both hands, plus additional parameters such as the translation between the hands and hand shape parameters.

We detail below the keypoint detection and encoding step, how we use the Keypoint Transformer to predict the hands poses, the representations we considered for the 3D hand poses, and the loss used for training. We also explain how our approach can be extended to object pose estimation during hand-object interaction scenarios.
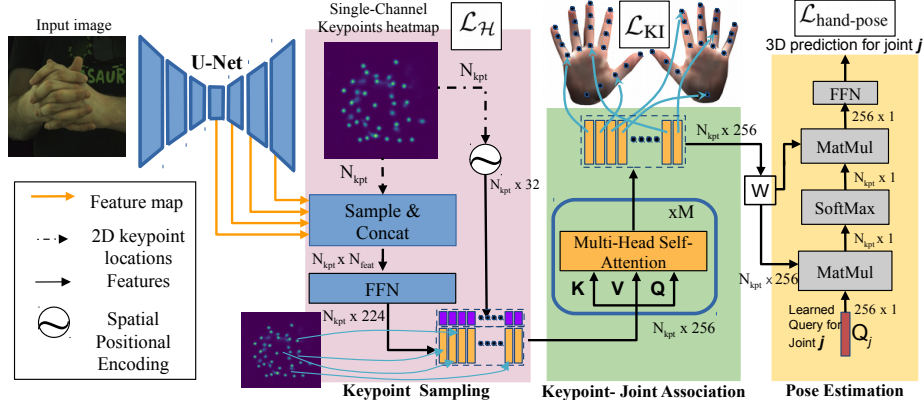
Figure 3. **Overview of our approach.** We detect keypoints which are potential locations of joints and encode them with CNN image features and spatial embedding (Section 3.1). From this information, the self-attention module creates context-aware keypoint features which are essential for associating each keypoint with the corresponding joint (Section 3.2). A cross-attention module finally predicts for each joint (using learned queries) the values required for computing the hand poses (Section 3.3). The exact nature of these values depends on the chosen representation of the hand pose (Section A). Not all the keypoints have to correspond to a joint and not all the joints have to be detected as keypoints, which makes our approach very robust but still accurate as it relies on keypoints (as discussed in Section 2).

## 3.1. Keypoint Detection and Encoding

Given the input image, we first extract keypoints that are likely to correspond to 2D hand joint locations. To do this, we predict a single-channel heatmap $\mathcal{H}$ from the input image using a standard U-Net [50] architecture, and we keep its local maximums using a non-differentiable, non-maximum suppression operation. At this stage, we do not attempt to recognize which keypoint corresponds to which joint as it is a difficult task, and the predicted heatmap has only one channel. In practice, we keep a maximum of $N_{\text{hand}}$ keypoints, with $N_{\text{hand}} = 64$, while the number of hand joints is 42 in total for 2 hands. The 2D keypoint locations are normalized to $[0, 1]$ range.

The ground truth heatmap $\mathcal{H}^*$ is obtained by applying a 2D Gaussian kernel of variance $\sigma$ at each of the ground truth 2D joint locations and the U-Net is trained to predict the heatmap by minimizing the L2 loss.

We compute for each detected keypoint an appearance and spatial encoding to represent the keypoints as input to the next stage. As shown in Fig. 3, for the appearance part, we extract image features from the decoder of the U-Net network. More exactly, we sample feature maps at multiple layers of the U-Net decoder at the normalized keypoint locations using bilinear interpolation and concatenate them to form a 3968-D feature vector, which is then reduced down to a 224-D encoding vector using a 3-layer MLP. For the spatial encoding, we obtain 32-D sine positional encoding similar to [8] corresponding to the 2D location of the keypoint. We finally concatenate the appearance and spatial encodings to form a 256-D vector representation of the keypoint. The keypoint detector is pre-trained before finetuning it jointly with the rest of the pipeline.

## 3.2. Keypoint-Joint Association

For each keypoint $K_i$, we have now an encoding vector $\mathcal{F}_i \in \mathbb{R}^{256}$. We use these vectors as input to the multilayer, multi-head self-attention module with $N_{SA}$ layers. The self-attention [61] helps to model the relationship between the keypoints and create global context-aware feature $\mathcal{G}_i \in \mathbb{R}^{256}$, for each keypoint. Such context-aware features are necessary to associate the keypoints with different joint types using a "keypoint-joint association" loss we denote $\mathcal{L}_{\text{KI}}$. As a result of $\mathcal{L}_{\text{KI}}$, the keypoint features also now encode the joint identity information along with the localized CNN image features.

The identity of keypoint $k$ is defined by $(h_k, j_k)$, where $h_k$ is the hand identity (left or right) and $j_k$ is the joint index. We also use an additional 'background' identity for keypoints that are falsely detected. The keypoint identity is predicted using a feed-forward network (FFN) consisting of a 2-layer MLP, a linear projection layer and a softmax layer. We use the cross-entropy loss for $\mathcal{L}_{\text{KI}}$:

$$\mathcal{L}_{\text{KI}} = \sum_i \text{CE}((h_i, j_i), (h_i^*, j_i^*)), \tag{1}$$

where $(h_i^*, j_i^*)$ are the ground truth identities and CE denotes the cross-entropy loss. To obtain the ground truth identity for the detected keypoints, we associate them at training time with the closest reprojection of a ground truth 3D joint, if the distance is below a threshold $\gamma$. If there are no reprojected joints within a distance of $\gamma$, the keypoint is assigned to the background class. We empirically set $\gamma = 3$ pixels in our experiments. Similar to [8], the keypoint identities are predicted after each layer of self-attention module using FFNs with shared weights and the loss is applied to predictions of each layer.

The prediction can result in multiple keypoints assigned to the same joint identity and some keypoints assigned to the background class. As we discuss in Section 5, the keypoints associated to the background are ignored, while all the keypoints associated with a given joint are considered for estimating the pose of the corresponding joint by the cross-attention module.

## 3.3. Pose Estimation from Identity-Aware Keypoints

The keypoint-joint association loss enables the keypoint features to also encode joint identity information along with the image features and spatial embeddings. We use a single cross-attention layer with learned joint queries to predict which keypoint(s) match the queried joint identities. The cross-attention operation [61] for a learned joint query $Q_j \in \mathbb{R}^{256}$ and features $\{\mathcal{G}_i\}_i$ is given by

$$\text{CA}(Q_j, G) = \text{softmax}\left(\frac{Q_j^T W_K G}{16}\right)(W_V G)^T, \quad (2)$$

where $G$ is a matrix whose columns contain feature vectors $\{\mathcal{G}_i\}_i$, and $W_K$ and $W_V$ are learnable matrices of dimension $256 \times 256$. Similar to [61], the cross-attention features are added to $Q_j$ to create a residual connection. The resulting features are transformed by a 3-layer MLP to map them to the pose space.

The number of joint queries depend on the pose representation. We consider 3 different representations and describe them in Section A and the suppl. mat. For example, we use 21 joint queries for each of the 21 joints per hand when using 2.5D pose representation. Along with the joint queries, one for each joint of the two hands, we use an additional learned query to predict the relative translation $T_{L \to R}$ between the hands, and the 10-D MANO hand shape parameters $\beta$. These are learned using the L1 loss. The MANO shape parameters are useful when predicting the pose using the MANO joint angle representation.

## 3.4. Hand Pose Representations and Losses

We consider three main hand pose representations: 3D joint locations, 2.5D [21, 37] joint locations, and MANO joint angles [49]. Previous methods [16,17,22,46,48] have noted that regressing model parameters such as joint angles is less accurate in terms of joint error than regressing the joint locations directly. However, regressing MANO joint angles provides access to the complete hand mesh required for modeling contacts and interpenetration during interactions [5, 17, 54] or for learning in a weakly supervised setup [3,16,27], which could be interesting for future extension of our method. As we show later in our experiments, the Keypoint Transformer enables the MANO joint angle representation to achieve competitive performance when
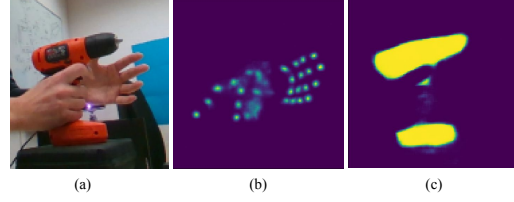

(a)        (b)        (c)

Figure 4. **Keypoint detection for hands and object.** We train a U-net decoder to predict (b) a heatmap for all the joints together and (c) a segmentation map for the object from each we extract keypoints at random locations.

compared to the joint location representation. We follow standard practice for these three representations. We detail them and their corresponding losses in the supplementary material for completeness.

## 3.5. Object Pose Estimation

Our method generalizes easily to predict the 3D pose of an object together with 3D poses of hands. As shown in Fig. 4, along with the heatmap for the hand keypoints, we also predict a segmentation map of the object by adding an additional prediction head to the U-Net decoder. We then randomly select $N_{\text{obj}} = 20$ points from this segmentation map and refer to them as 'object keypoints'. We also tried estimating the heatmap of 2D reprojections of fixed points on the object mesh and selecting their local maximums as object keypoints and obtained similar results.

We encode the appearance and spatial locations of the object keypoints in a 256-D vector, exactly like the hand keypoints. Collectively, these keypoint encodings cover the object appearance, allowing us to predict the 3D rotation and translation of the object. The encodings of $N_{\text{obj}}$ object keypoints and $N_{\text{hand}}$ hand keypoints are provided *together* to the self-attention module.

Along with the hand keypoint identities $(h_k, j_k)$ and the background identity described in Section 3.2, we rely on an additional identity for the object. During the keypoint association stage, all the keypoints originating from the object are associated with the 'object' identity, allowing the cross-attention module to only attend to object keypoints when estimating the object pose. Along with the joint queries that estimate the hand pose, we consider 2 additional queries in the cross-attention module and predict the 3D object rotation and 3D object translation relative to the right hand in a manner similar to that of hand pose. The object rotation is parameterized using the method proposed in [68].

We use a symmetry-aware object corner loss similar to [45] to train the network:

$$\mathcal{L}_{\text{obj-pose}} = \min_{R \in \mathcal{S}} \frac{1}{8} \sum_{i=1}^{8} ||P \cdot B_i - P^* \cdot R \cdot B_i||_2^2, \quad (3)$$

where $P$ and $P^*$ denote the estimated and ground-truth ob-

ject poses, $B_i$ denotes the $i^{\text{th}}$ corner of the 3D bounding box for the object in rest pose, and $\mathcal{S}$ is the set of rotation matrices which, when applied to the object, does not change its appearance.

## 3.6. End-to-End Training

We train our architecture end-end by minimizing the sum of the losses introduced above:

$$\mathcal{L} = \mathcal{L}_{\mathcal{H}} + \mathcal{L}_{\text{KI}} + \mathcal{L}_{\mathcal{T}} + \mathcal{L}_{\text{hand-pose}} + \mathcal{L}_{\text{obj-pose}} , \quad (4)$$

where $\mathcal{L}_{\text{hand-pose}}$ is the loss on the hand poses (detailed in the suppl. mat.) and $\mathcal{L}_{\mathcal{T}}$ is the L1 loss for relative translation between the two hands. Note that the keypoint detector is pretrained before training the entire network end-to-end. More optimization details are also given in the suppl. mat.

## 4. Evaluation

We evaluated our method on three challenging hand interaction datasets: InterHand2.6M, HO-3D, and our $H_2O$-3D dataset we introduce with this paper. We discuss them below.

### 4.1. InterHand2.6M

**Training and test sets.** InterHand2.6M [37] is a recently published two-hand interaction dataset with many challenging poses. It was annotated semi-automatically and contains 1.36M train images and 849K test images.

**Metrics.** We report the Mean Per Joint Position Error (MPJPE) and the Mean Relative-Root Position Error (MRRPE) to evaluate the root-relative hand pose and the translation between the hands respectively, as in [37].

**Baselines.** We consider two Transformer-based baseline architectures. The first baseline ('CNN+SA') provides the low-resolution ($32\times$ downsampled) CNN feature maps after flattening along the spatial dimensions as input to the Transformer encoder containing self-attention (SA) modules. The output tokens of the encoder are concatenated and the pose is predicted using an MLP. The second baseline ('CNN+SA+CA') is more similar to DETR [8], where

|  | MPJPE (mm) | | | MRRPE |
|  | Single Hand | Two Hands | All | (mm) |
| --- | --- | --- | --- | --- |
| CNN+SA | 13.53 | 16.87 | 15.31 | 33.84 |
| CNN+SA+CA (DETR [8]) | 12.81 | 15.94 | 14.48 | 32.87 |
| InterNet [37] | 12.16 | 16.02 | 14.22 | 32.57 |
| Ours | **10.99** | **14.34** | **12.78** | **29.63** |
| Dong et al. [25] | - | - | 12.08 | - |
| Ours | 9.10 | 11.98 | **11.30** | 21.89 |
| Fan et al. [12] | 11.32 | 15.57 | - | **30.51** |
| Ours | **11.08** | **15.33** | 13.41 | 30.87 |

Table 1. **Comparison with 2 baselines and the state-of-the-art methods on InterHand2.6M [37].** We compare with [12, 25, 37] using the different train/test splits reported in their works.

the low-resolution CNN feature maps are provided to the Transformer encoder-decoder architecture. The Transformer decoder contains SA and cross-attention (CA) modules. The queries in the decoder are learnt and the pose is predicted using FFN similar to our Keypoint Transformer. We provide more details about the baselines in the suppl. mat. These baselines help to understand the importance of keypoint sampling and selection for pose estimation.

**Results.** Table 1 compares the accuracy of our method with the state-of-the-art method InterNet [37], and the two baselines, when using the 2.5D pose representation. Our method achieves 10% higher accuracy than InterNet, which is a CNN-based architecture, and 16% and 12% higher accuracy than the two baselines, respectively. The higher accuracy of 'CNN+SA+CA' w.r.t 'CNN+SA' baseline demonstrates that soft-selection of image features by the decoder improves the accuracy. Further, the higher accuracy (12%) of our Keypoint-Transformer w.r.t the 'CNN+SA+CA' architecture shows that use of keypoint features for pose estimation instead of features from the entire image increases the overall accuracy.

We compare our method with [25] and [12] using their train and test splits. [12, 25] use per-joint heatmaps coupled with joint visibility and segmentation guided features to improve the accuracy of the pose estimation, thus resulting in the model complexity that is higher than InterNet [37]. Our method with a model complexity same as [37] (see Section 5) still outperforms these state-of-the-art methods. We show qualitative results in Fig. 5 and in the suppl. mat.

### 4.2. HO-3D

**Training and test sets.** The HO-3D [14] dataset contains automatically annotated hand-object interaction sequences of a right hand and an object from the YCB [64] dataset. It contains 66K training images and 11K test images. We consider only objects seen in the training set for evaluation.

**Metrics.** As in [14], we report the mean joint error after scale-translation alignment of the root joint and the area-under-the-curve (AUC) metrics to evaluate the hand pose. The object pose is computed w.r.t to the hand frame of reference and is evaluated using the standard Maximum Symmetry-Aware Surface Distance (MSSD) [19], as it considers the symmetricity of objects.

**Results.** We use 3D joint representation to estimate the hand pose. Table 2 compares the accuracy of the proposed hand-object pose estimation method with other approaches. Keypoint Transformer performs significantly better than previous methods [14, 16, 17]. As [14, 16, 17] do not consider symmetricity of objects during training and evaluation, we also report our results in a similar setting. We show qualitative results in Fig. 6. Please refer to supplementary material for quantitative comparison with more recent works.

Figure 5. **Qualitative results on InterHand2.6M [37].** Our method obtains accurate poses of hands during complex interactions. We show the estimated MANO model from a different view.



HO-3D                    H$_2$O-3D

Figure 6. **Qualitative results for our method on the H$_2$O-3D and HO-3D datasets**. Our method recovers poses even under large occlusions by the object and achieves state-of-the-art results on HO-3D while serving as a strong baseline for our new dataset H$_2$O-3D. Note that some objects (columns 2&4) are considered to be rotationally symmetric along the z-axis.

|      | Camera Intrinsics | Image Crop | Joint Error | Mean Joint AUC | MSSD (Object Pose Error) |
|------|-------------------|------------|-------------|----------------|--------------------------|
| [14] | Yes               | Yes        | 3.04        | 0.49           | -                        |
| [17] | No                | Yes        | 3.18        | 0.46           | -                        |
| [16] | Yes               | No         | 3.69        | 0.37           | 11.99                    |
| Ours | No                | Yes        | **2.57**    | **0.54**       | **7.02**                 |

Table 2. **Accuracy of our method on the HO-3D dataset for hand and object pose estimation.** Our method outperforms previous methods by a large margin.

### 4.3. H$_2$O-3D

**Training and test sets.** We introduce a dataset named H$_2$O-3D comprising sequences of two hands manipulating an object automatically annotated with the 3D poses of the hands and the object, by extending the work of [14] to consider two hands. Figs. 1 and 6 show some images. Five different subjects manipulate 10 different objects from YCB using both hands with a functional intent. We captured 60'998 training images and 15'342 test images using a 5 RGBD cameras multi-view setup. The H$_2$O-3D test set contains 7 objects seen in the training set and 1 unseen object. More details are provided in the supplementary material. H$_2$O-3D is significantly more challenging than previous hand interaction datasets as there are many large occlusions between the hands and the objects.

**Metrics and Results** We use the 3D joint representation for the hand pose and evaluate the accuracy using the MPJPE and MRRPE metrics (see Section 4.1) for the hand and the MSSD metric for the object (see Section 4.2). Details about the angle of symmetry for different objects considered during training and evaluation is provided in the suppl. mat. Due to large occlusions of the object by the hands, a portion of images are unsuitable for object pose estimation. We identify these images as the ones whose ground truth object segmentation area is less than 2% of the cropped image area and exclude them from the object pose estimation during training and evaluation. We also used the HO-3D train split and mirrored the images randomly during training to obtain right hand- and left hand-only images, to later combine with the training set of H$_2$O-3D.

Our method achieves an MPJPE of 3.09 cm and an MR-RPE of 8.28 cm on this dataset. Due to large occlusions by the object, estimating the translation between the hands is more challenging and the MRRPE is about 2.5 times worse than on InterHands2.6M which does not contain objects. On objects, our method achieves MSSD values of 7.96 cm. We provide object-specific MSSD values in the supplementary material. Fig. 6 shows qualitative results.
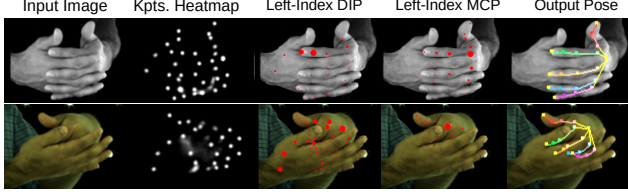
Figure 7. **Visualizing the cross-attention weights for two joint queries of the left hand.** The radius of the red circles are proportional to the weights. When the joint is occluded like the DIP joint on the second row, nearby visible keypoints are selected by the attention mechanism for pose estimation.
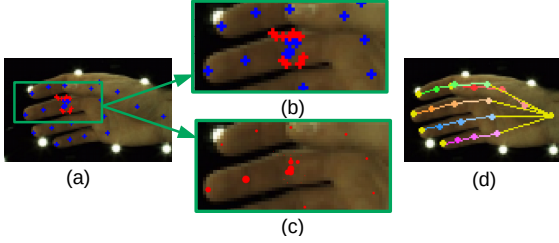


Figure 8. **Robustness to noisy keypoints.** In this example, we added noisy keypoints around the middle finger PIP joint. Most of the noisy keypoints are predicted to belong to background class (in red), while some are associated with the PIP joint (in blue). The noisy keypoints associated with the PIP joint have all higher cross-attention weights (c) and are considered for final pose estimation.
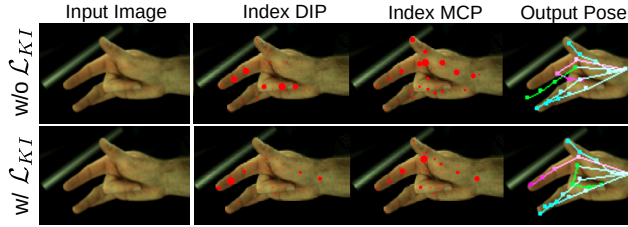


Figure 9. **Cross-attention with and without the keypoint-joint association loss $\mathcal{L}_{KI}$.** $\mathcal{L}_{KI}$ makes the keypoints 'identity-aware', resulting in higher accuracy.

## 5. Discussion

We report here the results of experiments we perform using InterHand2.6M (V0.0) to understand better our method.

**Visualization of cross-attention.** We visualize the cross-attention weights in Fig. 7 for two joint queries of the left-index finger. When the joint is not occluded, as in the first row, each joint query attends to the keypoint whose location coincides with the corresponding joint location in the image. In other words, local image features at the location of the joint are used to estimate the pose of that joint. We believe this property of using local image features helps in achieving higher accuracy than other CNN-based approaches [16, 17, 37]. In the second row of Fig. 7, the left index finger is occluded except for the MCP joint and no keypoints are detected for the invisible joints. The cross-attention module *selects* nearby visible keypoints resulting

|  | Camera Intrinsics | MPJPE (mm) | | | MRRPE (mm) |
|---|---|---|---|---|---|
|  |  | Single Hand | Two Hands | All |  |
| 3D | No | 12.42 | **17.08** | 14.76 | **33.14** |
| 2.5D | Yes | **11.73** | 17.69 | **14.73** | 34.40 |
| $\theta$ | No | 15.36 | 20.61 | 18.01 | 37.91 |

Table 3. **Accuracy obtained with the 3 different pose representations.**

|  | $N_{CA} = 1$, Varying $N_{SA}$ | | | $N_{SA} = 6$, Varying $N_{CA}$ | | | [37] |
|---|---|---|---|---|---|---|---|
|  | 0 | 3 | 6 | 1 | 3 | 6 |  |
| Single Hand | 12.34 | 11.77 | 11.24 | 11.24 | 11.14 | 11.08 | 12.63 |
| Two Hands | 16.93 | 15.55 | 15.44 | 15.44 | 15.35 | 15.33 | 17.36 |

Table 4. **3D pose accuracy (MPJPE, in mm) for different numbers of self-attention ($N_{SA}$) and cross-attention ($N_{CA}$) layers.**

|  | Resnet-18 | Resnet-34 | Resnet-50 | [37] (Resnet-50) |
|---|---|---|---|---|
| Total Params | 28M | 38M | 48M | 48M |
| Single Hand | 11.67 | 11.99 | 11.28 | 12.63 |
| Two Hands | 16.78 | 16.41 | 15.32 | 17.36 |

Table 5. **3D pose accuracy (MPJPE, in mm) for different backbones.**

in a more global-level feature for estimating the joint pose.

**Robustness to noisy keypoints.** To demonstrate the robustness of our method, we added incorrect keypoints around the detected keypoints. As shown in Fig. 8 and supplementary material, most of these keypoints are labeled as background and all the keypoints that are assigned to the same joint are considered equally for the pose estimation.

**Importance of the keypoint-joint association loss $\mathcal{L}_{KI}$.** $\mathcal{L}_{KI}$ helps the cross-attention module to select the appropriate features for pose estimation as visualized in Fig. 9. Further, $\mathcal{L}_{KI}$ improves MPJPE by 10% (17.08 mm v/s 18.91 mm) and MRRPE by 15% (33.14 mm v/s 38.96 mm) on interacting hand images.

**Accuracy with different pose representations.** Table 3 compares the accuracy of the 3 hand pose representations that we consider. While the accuracy of the 3D and 2.5D representations are similar, the joint angle representation results in lower accuracy, in line with the observation from previous works [16, 17, 22, 46, 48].

**Effect of the number of self-attention (SA) and cross-attention (CA) layers.** Table 4 reports the MPJPE with different combinations of SA and CA layers. Even in the absence of any SA layers, our method outperforms [37]. Adding more CA layers has little effect on the accuracy.

**Effect of the number of parameters.** Table 5 reports the MPJPE for different CNN backbones. While larger backbones improve the accuracy, our method outperforms [37] even with a Resnet-18 backbone with approximately half the total number of parameters.

## 6. Conclusion

We showed that, by integrating a keypoint detector into a Transformer architecture, we could predict 3D poses of hands and objects from very challenging images, in a much

more accurate way than a standard Transformer architecture does. As we rely on keypoints, we believe that our approach is more general and could be applied to other problems, such as human and other articulated objects pose prediction and object category pose prediction [47].

# Supplementary Material

In this supplementary material, we discuss the limitations of our method, provide more details about the experiments and also show several qualitative results and comparisons. We also refer the reader to the **Supplementary Video** for visualization of results on different action sequences.

## A. Hand Pose Representations and Losses

We detail the three possible representations mentioned in Section 3.4 of the paper. We assume 21 3D-joint locations per hand as in the MANO [49] model. The losses for each of the 3 representations are summarized in Table 6.

**3D representation.** In this representation, each joint $j$ is associated with a parent-relative joint vector $V(j) = J_{3D}(j) - J_{3D}(p(j))$, where $J_{3D}$ is the 3D joint location and $p(j)$ refers to the parent joint index of joint $j$. We estimate 20 joint vectors per hand using 20 joint queries, one for each skeletal bone (40 queries for two hands), from which we can compute the root-relative 3D location, $J_{3D}^r$ of each joint by simple accumulation. The advantage of this representation is that it defines the hand pose relative to its root without requiring knowledge of the camera intrinsics.

**2.5D representation [21, 37].** In this representation, each joint is parameterised by its 2D location $J_{2D}$, and the difference $\Delta Z^p$ between its depth and the depth of its parent joint. The camera intrinsics matrix $K$ and the absolute depth $Z_{\text{root}}$ of the root joint (the wrist) [37] or the scale of the hand [21] are then required to reconstruct the 3D pose of the hand in camera coordinate system as $J_{3D} = K^{-1} \cdot (Z_{\text{root}} + \Delta Z^r) \cdot \left[ J_{2D_x}, J_{2D_y}, 1 \right]^T$, where $\Delta Z^r$ is the root-relative depth of the joint computed from its predicted $\Delta Z^p$ and the predicted $\Delta Z^p$ for its parents. $J_{2D_x}, J_{2D_y}$ are the predicted $x$ and $y$ coordinates of $J_{2D}$.

When using this representation, we also predict the root depth $Z_{\text{root}}$ separately using RootNet [35] as in [37]. Each joint query estimates the $J_{2D}$ and $\Delta Z^r$ for that joint and we require a total of 21 joint queries (42 for two hands), one for each joint location to estimate the 2.5D pose per hand.

**MANO joint angles, $\theta$ [49].** In this representation, each 3D hand pose is represented by 16 3D joint angles in the hand kinematic tree and is estimated using 16 joint queries

| Representation | $\mathcal{L}_{hand-pose}$ |
|---|---|
| 3D | $\sum_j \lVert V(j) - V(j)^* \rVert_1 + \sum_j \lVert J_{3D}^r(j) - J_{3D}^{r^*}(j) \rVert_1$ |
| 2.5D | $\sum_j \lVert J_{2D}(j) - J_{2D}^*(j) \rVert_1 + \sum_j \lvert \Delta Z^r(j) - \Delta Z^{r^*}(j) \rvert$ |
| $\theta$ | $\sum_j \lVert J_{3D}^r(j) - J_{3D}^{r^*}(j) \rVert_1 + \sum_j \lVert \theta(j) - \theta^*(j) \rVert_1$ |

Table 6. Hand pose losses for different pose representations. $x^*$ denotes the ground-truth values for variable $x$ and $x(j)$ the value of $x$ at joint $j$.

per hand, one for each joint. The MANO hand shape parameter is estimated along with the relative translation between the hands using an additional query. Given the predicted 3D joint angles $\theta$ for each hand and the shape parameters $\beta$, it is possible to compute the root-relative 3D joint locations, $J_{3D}^r$ of each hand.

## B. Comparison with state-of-the-art on HO3D Dataset

We compare the performance our method with several other methods on the HO-3D(V2) and HO-3D(V3) datasets and show the results in Tab. 7 and Tab. 8, respectively. On HO-3D(V2), the performance of our method is very close to the HandOccNet [44], which achieves the highest accuracy when considering the scale-translation aligned MPJPE metric. However, HandOccNet achieves higher accuracy when considering the procrustes aligned MPJPE metric. On the HO-3D(V3) dataset, our method performs worse than Arti-Boost [31], which uses additional training data.

## C. Method Limitations

Though our method results in accurate poses during interactions, the results are sometimes not plausible as we do not model contacts and interpenetration [5, 17, 23] between hands and objects. Further, during highly complex and severely occluded hand interactions as we show in the last row of Fig. 18, our method fails to obtain reasonable hand poses. We believe these problems can be tackled in the future by incorporating temporal information and physical modeling into our architecture.

## D. Hand-Object Pose Estimation Pipeline

In Fig. 10, we show the complete pipeline of our Keypoint-Transformer architecture for estimating poses of two hands and object during interaction.

## E. Implementation details

The encoder of our U-Net [50] is based on ResNet-50 [18] architecture while a series of upsampling and convolutional layers with skip connections forms the U-Net decoder. We use 256×256 pixels as input image resolution, 128×128 pixels as heatmap resolution, and set the 2D Gaussian kernel variance, $\sigma$ to 1.25 during training. The

| Method | Pose Repr. | Joint Error (scale and trans. align.) in cms | Joint Error AUC (scale and trans. align.) | Joint Error (Procrustes align.) in cms | Joint Error AUC (Procrustes align.) | Mesh Error (Procrustes align.) in cms | Mesh Error AUC (Procrustes align.) | Mesh Error F-Score @5mm | Mesh Error F-Score @15mm |
|---|---|---|---|---|---|---|---|---|---|
| METRO [32] | Mesh | 2.89 | 0.504 | 1.04 | 0.792 | 1.11 | 0.779 | 0.484 | 0.946 |
| Liu et al. [34] | Joint angle | 3.17 | 0.463 | 0.99 | 0.803 | 0.95 | 0.810 | 0.528 | 0.956 |
| HandOccNet [44] | Joint angle | 2.40 | 0.557 | 0.91 | 0.819 | 0.88 | 0.819 | 0.564 | 0.963 |
| I2L-MeshNet [36] | Mesh | 2.60 | 0.529 | 1.12 | 0.775 | 1.39 | 0.722 | 0.409 | 0.932 |
| Pose2Mesh [10] | Mesh | 3.33 | 0.480 | 1.25 | 0.754 | 1.27 | 0.749 | 0.441 | 0.909 |
| I2UV-HandNet [9] | Mesh | - | - | 0.99 | 0.804 | 1.01 | 0.799 | 0.500 | 0.943 |
| Zheng et al. [67] | 2.5D | 2.51 | 0.541 | - | - | - | - | - | - |
| Hampali et al. [14] | 3D | 3.04 | 0.494 | 1.07 | 0.788 | 1.06 | 0.790 | 0.506 | 0.942 |
| Hasson et al. [17] | Joint angle | 3.18 | 0.461 | 1.10 | 0.780 | 1.12 | 0.777 | 0.464 | 0.939 |
| Hasson et al. [16] | Joint angle | 3.69 | 0.369 | 1.14 | 0.773 | 1.14 | 0.773 | 0.428 | 0.932 |
| ArtiBoost [31] | 2.5D | 2.53 | 0.532 | 1.14 | 0.773 | 1.09 | 0.782 | 0.488 | 0.944 |
| Ours-ResNet18 | Joint angle | 3.11 | 0.459 | 1.10 | 0.780 | 1.13 | 0.774 | 0.444 | 0.935 |
| Ours-ResNet50 | 3D | 2.57 | 0.553 | 1.08 | 0.786 | - | - | - | - |

Table 7. Comparison with state-of-the-art methods on HO-3D V2 dataset.

| Method | Pose Repr. | Joint Error (scale and trans. align.) in cms | Joint Error AUC (scale and trans. align.) | Joint Error (Procrustes align.) in cms | Joint Error AUC (Procrustes align.) | Mesh Error (Procrustes align.) in cms | Mesh Error AUC (Procrustes align.) | Mesh Error F-Score @5mm | Mesh Error F-Score @15mm |
|---|---|---|---|---|---|---|---|---|---|
| ArtiBoost [31] | 2.5D | 2.34 | 0.565 | 1.08 | 0.785 | 1.04 | 0.792 | 0.507 | 0.946 |
| Ours-ResNet50 | 3D | 2.48 | 0.575 | 1.09 | 0.785 | - | - | - | - |

Table 8. Comparison with state-of-the-art methods on HO-3D V3 dataset. Note that ArtiBoost [31] uses additional training data which is not used in our method.

$256 \times 256$ pixel input image patch is loosely cropped around the hand and object. We use Adam [26] optimizer with a learning rate of $10^{-4}$ and $10^{-5}$ for the attention modules and CNN backbone, respectively. The network is trained for 50 epochs on 3 Titan V GPUs with a total batch size of 78 and uses on-line augmentation techniques such as rotation, scale and mirroring during training.

## F. Baseline Architectures

We detail here the two baselines, 'CNN+SA' and 'CNN+SA+CA' considered in Section 4.1 of the main paper. Figures 11 and 12 show their architectures. We used $256 \times 256$ cropped images as input to the CNN resulting in a feature map of spatial dimensions $8 \times 8$ and 2048 channels. The features are flattened along the spatial dimensions and the 64 features are converted to 224 dimensions using 3 MLP layers. These features are then concatenated with 32-D positional embeddings resulting in 256-D features and are provided to the Transformer encoder. The networks were trained to output the 2.5D pose representation for 50 epochs on 3 Titan V GPUs with a batch size of 78. The joint queries in 'CNN+SA+CA' are learned in a similar way as for our Keypoint Transformer.

## G. Robustness to Noisy Keypoints

We show more examples to demonstrate the robustness of our method to noisy keypoints. We consider two scenarios, adding noisy keypoints to the set of detected keypoints, and randomly removing some keypoints from the set of detected keypoints. We show results in Figures 13 and 14, respectively. The number of detected keypoints for these cases were 48 and we added 30 additional noisy keypoints for the former scenario and retained only 30 keypoints for the latter scenario.

## H. H$_2$O-3D Dataset

Our dataset contains sequences of two hands interacting with an object, captured on a multi-view setup with 5 RGBD cameras. We collected data from six different subjects and considered ten objects from the YCB dataset with each subject manipulating the object with a functional intent. The dataset is automatically annotated with 3D poses of hands and objects using the optimization method of [14]. The dataset contains 60'998 training images and 15'342 test images from 17 different multi-view sequences in total. As explained in the main paper, we only consider 9'098 images from the set of 15'342 test images for object pose evaluation as the objects in the remaining images are barely visible due to occlusion by the hands. We show some sample annotations from the dataset in Fig. 15. Table 9 shows the list of YCB objects and their axis and angle of symmetry considered during our training and evaluation.

### H.1. Per-Object MSSD Values with Keypoint Transformer

Table 10 shows the accuracy of the object poses estimated by our Keypoint Transformer on the H$_2$O-3D dataset using the MSSD metric as described in Section 4.3 of the main paper.
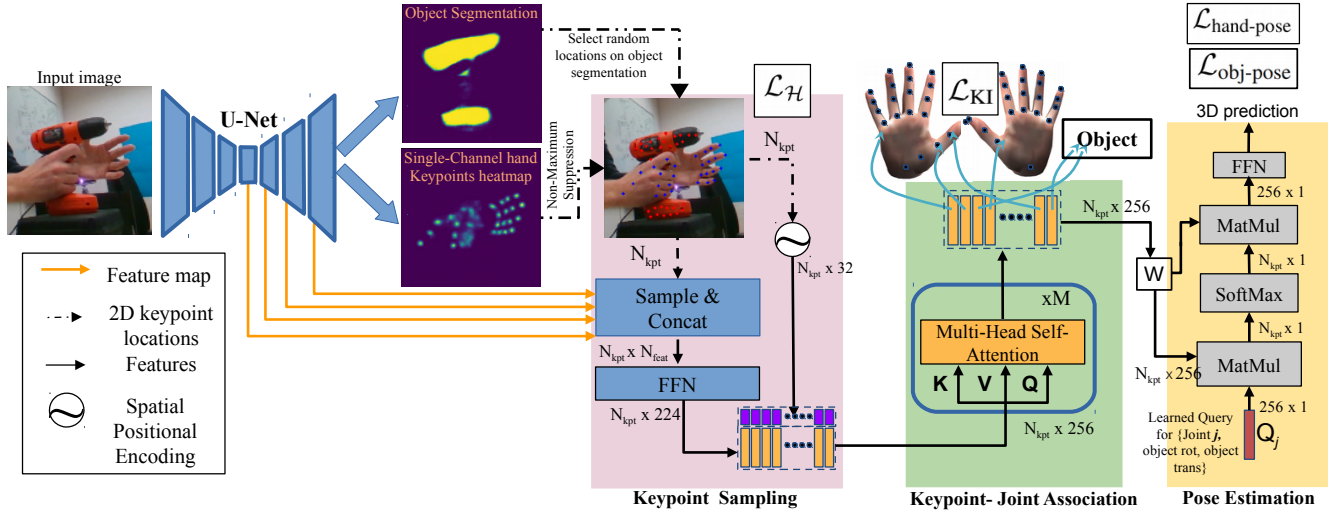
Figure 10. **Pipeline for hands and object pose estimation.** The object keypoints are selected by randomly sampling 2D locations on the object segmentation map regressed by the U-Net (Section 3.5 of main paper). The hand keypoints are selected from the single-channel keypoints heatmap, also regressed by the U-Net (Section 3.1 of main paper). Each of the detected keypoints are encoded using CNN image features and spatial embedding. The keypoints are associated with one of the 42 hand joints (21 joints per hand), the object class or the background class in the keypoint-joint association stage (Section 3.2 of main paper). The object rotation and translation w.r.t the right hand is estimated in the pose estimation stage using 2 different learned object queries, while the pose of each hand-joint is estimated using per-joint learned queries (Section 3.3 of main paper).
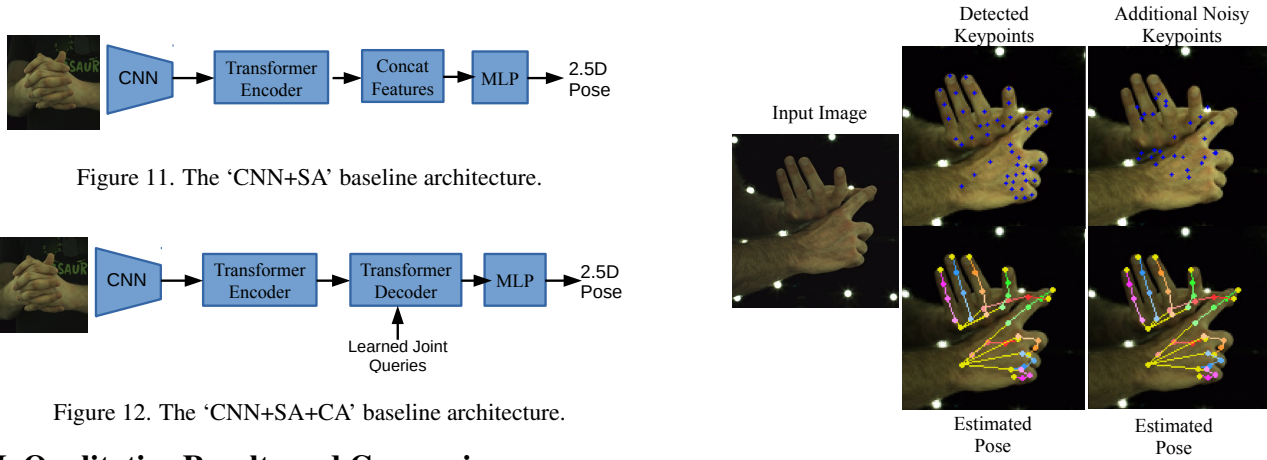


Figure 11. The 'CNN+SA' baseline architecture.



Figure 12. The 'CNN+SA+CA' baseline architecture.

## I. Qualitative Results and Comparisons

We provide here more qualitative results on HO-3D, $H_2$O-3D and InterHand2.6M.

### I.1. HO-3D and $H_2$O-3D Qualitative Results

Fig. 16 shows qualitative results on $H_2$O-3D and HO-3D. Note that as we do not model contacts and interpenetration between hands and object, our method sometimes results in implausible poses as we show in the last example of Fig. 16.

### I.2. InterHand2.6M Qualitative Results

Fig. 17 compares the estimated poses using the InterNet method from [37] and our proposed approach. As noted in Section 1 and Table 3 of the main paper, purely CNN-based



Figure 13. Effect of adding additional noisy keypoints. Our method predicts accurate poses even with noisy keypoints.

approaches do not explicitly model the relationship between image features of joints and tend to *confuse* joints during complex interactions. Our method performs well during complex interactions and strong occlusions (see last row of Fig. 17).

We show more qualitative results using the MANO angle representation in Fig. 18. Our retrieved poses are very similar to ground-truth poses. As we show in the last row of Fig. 18, our method fails during scenarios where the hand is severely occluded during complex interaction.
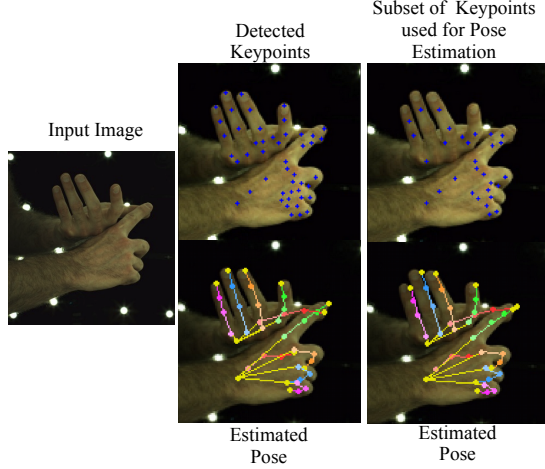
Figure 14. Effect of using a subset of detected keypoints for pose estimation. We consider only 30 of the 48 detected keypoints for pose estimation and still estimate an accurate pose.

## J. Attention Visualization

In Fig. 19, we show more visualization of the cross-attention weights for three different joint queries. More specifically, the cross-attention weights represent the multi-

| Object | Axis | Angle |
|---|---|---|
| Mustard Bottle | Z | $180^o$ |
| Bleach Cleanser | Z | $180^o$ |
| Cracker Box | Z | $180^o$ |
| Sugar Box | Z | $180^o$ |
| Potted Meat Can | Z | $180^o$ |
| Bowl | Z | $\infty$ |
| Mug | Z | $\infty$ |
| Pitcher Base | Z | $\infty$ |
| Banana | - | - |
| Power Drill | - | - |

Table 9. $H_2$O-3D objects and their axis and angle of symmetry considered during training and evaluation with our Keypoint Transformer.

| Object | MSSD (cm) |
|---|---|
| Bleach Cleanser | 7.7 |
| Mug | 6.5 |
| Banana | 9.8 |
| Pitcher Base | 7.9 |
| Bowl | 7.8 |
| Scissors | 13.5 |
| Power Drill | 8.5 |
| All | 7.9 |

Table 10. Object pose estimation accuracy of our Keypoint Transformer on the $H_2$O-3D dataset.

plicative factor on each of the keypoint features for a given joint query. We observe that the cross-attention learns to select keypoint(s) from respective joint location for each joint query when the joint is visible. For occluded joints, features from nearby visible joints are selected.

## References

[1] Anil Armagan, Guillermo Garcia-Hernando, Seungryul Baek, Shreyas Hampali, Mahdi Rad, Zhaohui Zhang, Shipeng Xie, Ming-xiu Chen, Boshen Zhang, F. Xiong, Yang Xiao, Zhiguo Cao, Junsong Yuan, Pengfei Ren, Weiting Huang, Haifeng Sun, Marek Hrúz, Jakub Kanis, Zdenek Krnoul, Qingfu Wan, Shile Li, Linlin Yang, Dongheui Lee, Angela Yao, Weiguo Zhou, Sijia Mei, Yunhui Liu, Adrian Spurr, Umar Iqbal, Pavlo Molchanov, Philippe Weinzaepfel, Romain Brégier, Grégory Rogez, Vincent Lepetit, and Tae-Kyun Kim. 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. 2

[2] Seungryul Baek, Kwang In Kim, and Tae-Kyun Kim. Pushing the Envelope for RGB-Based Dense 3D Hand Pose Estimation via Neural Rendering. In *Conference on Computer Vision and Pattern Recognition*, 2019. 1

[3] Seungryul Baek, Kwang In Kim, and Tae-Kyun Kim. Weakly-Supervised Domain Adaptation via GAN and Mesh Model for Estimating 3D Hand Poses Interacting Objects. In *Conference on Computer Vision and Pattern Recognition*, 2020. 5

[4] Luca Ballan, Aparna Taneja, Juergen Gall, Luc Van Gool, and Marc Pollefeys. Motion Capture of Hands in Action Using Discriminative Salient Points. In *European Conference on Computer Vision*, 2012. 2

[5] Samarth Brahmbhatt, Chengcheng Tang, Christopher D. Twigg, Charles C. Kemp, and James Hays. ContactPose: A Dataset of Grasps with Object Contact and Hand Pose. In *European Conference on Computer Vision*, 2020. 2, 3, 5, 9

[6] Guillem Braso, Nikita Kister, and Laura Leal-Taixé. The Center of Attention: Center-Keypoint Grouping Attention for Multi-Person Pose Estimation. In *International Conference on Computer Vision*, 2021. 3

[7] Zhe Cao, Ilija Radosavovic, Angjoo Kanazawa, and Jitendra Malik. Reconstructing Hand-Object Interactions in the Wild. In *arXiv Preprint*, 2020. 2, 3

[8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End Object Detection with Transformers. In *European Conference on Computer Vision*, 2020. 2, 3, 4, 6

[9] Ping Chen, Yujin Chen, Dong Yang, Fangyin Wu, Qin Li, Qingpei Xia, and Yong Tan. I2uv-handnet: Image-to-uv prediction network for accurate and high-fidelity 3d hand mesh modeling. In *International Conference on Computer Vision*, 2021. 10

[10] Hongsuk Choi, Gyeongsik Moon, and Kyoung Mu Lee. Pose2mesh: Graph convolutional network for 3d human pose and mesh recovery from a 2d human pose. In *European Conference on Computer Vision (ECCV)*, 2020. 10

Figure 15. Samples from H$_2$O-3D dataset. Our dataset contains sequences with complex actions performed by both hands on YCB [64] objects.



Figure 16. Qualitative results on H$_2$O-3D and HO-3D [14]. Our method obtains state-of-the-art results on HO-3D while predicting reasonable results on H$_2$O-3D. The last example is a failure case where the predicted relative translations are inaccurate.
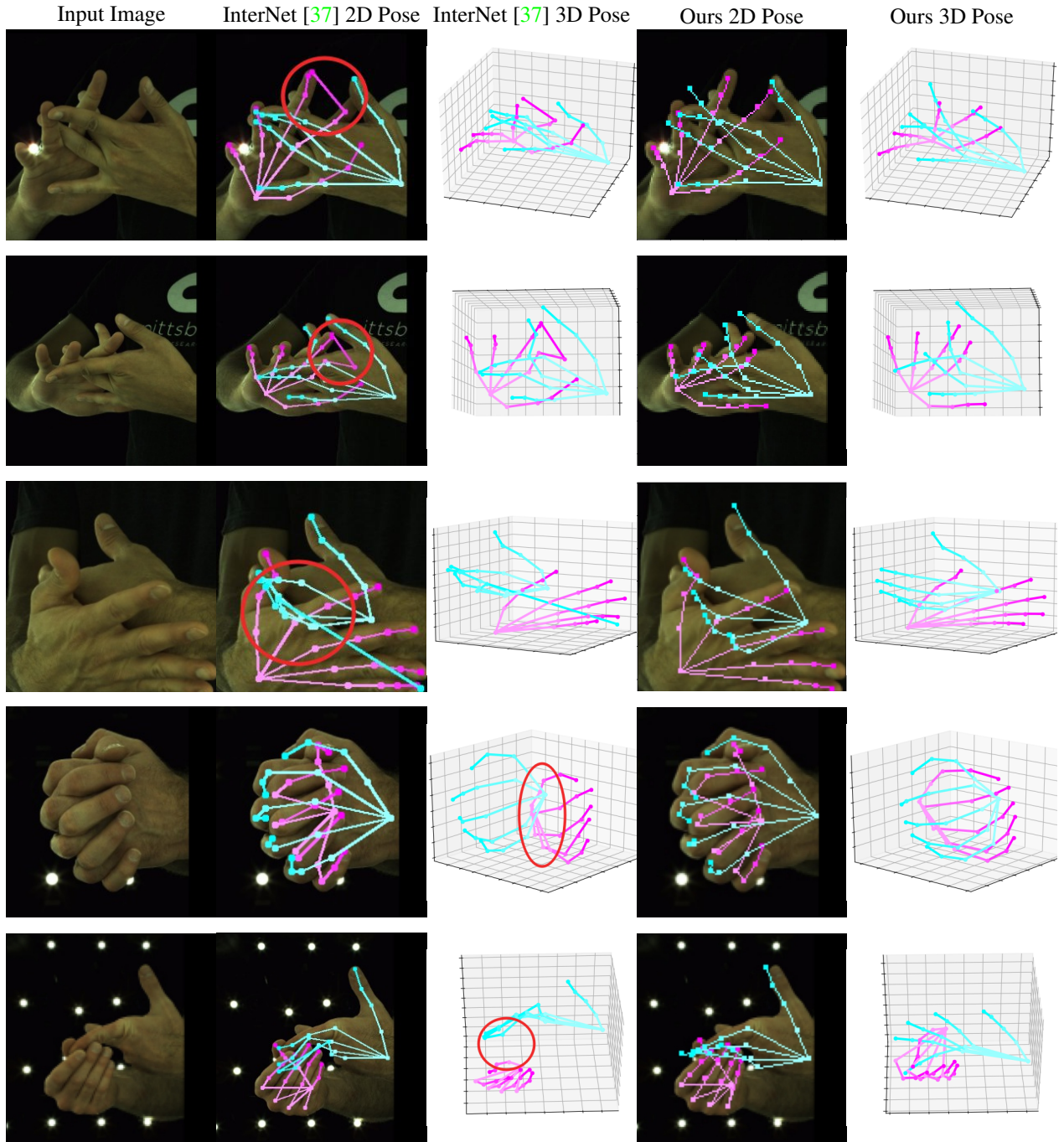
Figure 17. Qualitative comparison between InterNet [37] and our proposed method. Our method outputs more accurate poses even during strong occlusions. Red circles indicate regions where InterNet results are inaccurate.

[11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *arXiv Preprint*, 2020. 3

[12] Zicong Fan, Adrian Spurr, Muhammed Kocabas, Siyu Tang, Michael Black, and Otmar Hilliges. Learning to Disambiguate Strongly Interacting Hands via Probabilistic Per-Pixel Part Segmentation. In *International Conference on 3D Vision*, 2021. 2, 3, 6

[13] Guillermo Garcia-Hernando, Shanxin Yuan, Seungryul Baek, and Tae-Kyun Kim. First-Person Hand Action Benchmark with RGB-D Videos and 3D Hand Pose Annotations. In *Conference on Computer Vision and Pattern Recognition*, 2018. 3

Figure 18. Qualitative results of our method on InterHand2.6M [37] compared to ground-truth poses. Our method predicts accurate poses in most scenarios. The last row shows a failure case where our method cannot recover the accurate pose due to complex pose and severe occlusion.
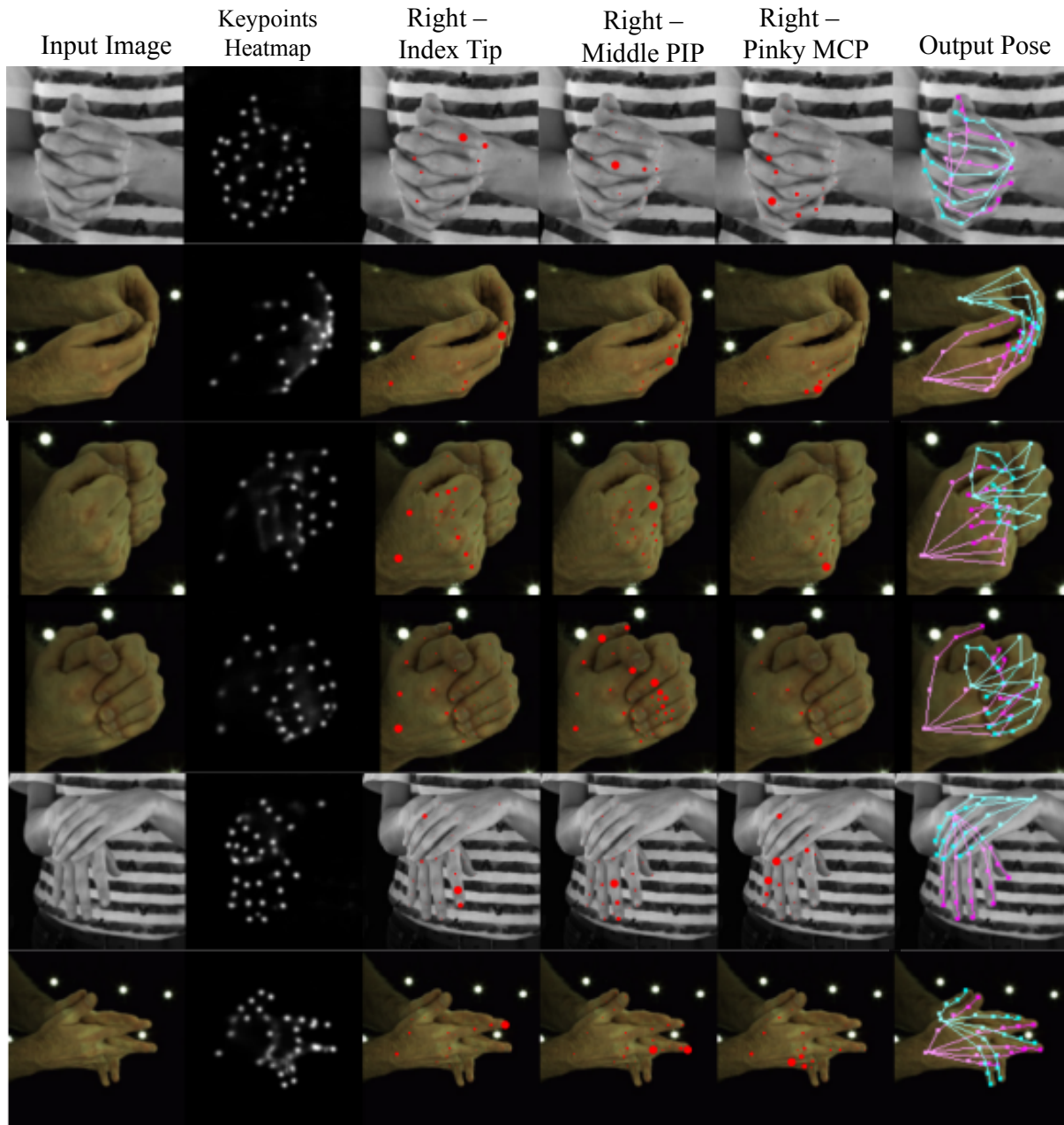
Figure 19. Attention visualization for 3 joint queries. Each joint query attends to the image feature from the respective joint location.

[14] Shreyas Hampali, Mahdi Rad, Markus Oberweger, and Vincent Lepetit. HOnnotate: A Method for 3D Annotation of Hand and Object Poses. In *Conference on Computer Vision and Pattern Recognition*, 2020. 1, 2, 3, 6, 7, 10, 13

[15] Shangchen Han, Beibei Liu, R. Cabezas, Christopher D. Twigg, Peizhao Zhang, Jeff Petkau, Tsz-Ho Yu, Chun-Jung Tai, Muzaffer Akbay, Zheng Wang, Asaf Nitzan, Gang Dong, Yuting Ye, Lingling Tao, Chengde Wan, and Robert Wang. MEgATrack: Monochrome Egocentric Articulated Hand-Tracking for Virtual Reality. *IEEE Transactions on Robotics and Automation*, 39, 2020. 1, 2, 3

[16] Yana Hasson, Bugra Tekin, Federica Bogo, Ivan Laptev, Marc Pollefeys, and Cordelia Schmid. Leveraging Photometric Consistency over Time for Sparsely Supervised Hand-Object Reconstruction. In *Conference on Computer Vision and Pattern Recognition*, 2020. 2, 3, 5, 6, 7, 8, 10

[17] Yana Hasson, Gül Varol, Dimitrios Tzionas, Igor Kalevatykh, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning Joint Reconstruction of Hands and Manipulated Objects. In *Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2, 3, 5, 6, 7, 8, 9, 10

[18] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Conference on Computer Vision and Pattern Recognition*, 2016. 9

[19] Tomás Hodan, Martin Sundermeyer, Bertram Drost, Yann Labbé, Eric Brachmann, Frank Michel, Carsten Rother, and Jiri Matas. BOP Challenge 2020 on 6D Object Localization. In *Computer Vision - ECCV 2020 Workshops - Glasgow, UK, August 23-28, 2020, Proceedings, Part II*, 2020. 6

[20] Lin Huang, Jianchao Tan, Ji Liu, and Junsong Yuan. Hand-Transformer: Non-Autoregressive Structured Modeling for 3D Hand Pose Estimation. In *European Conference on Computer Vision*, 2020. 3

[21] Umar Iqbal, Pavlo Molchanov, Thomas Breuel, Juergen Gall, and Jan Kautz. Hand Pose Estimation via Latent 2.5D Heatmap Regression. In *European Conference on Computer Vision*, 2018. 1, 5, 9

[22] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-End Recovery of Human Shape and Pose. In *Conference on Computer Vision and Pattern Recognition*, 2018. 5, 8

[23] Korrawe Karunratanakul, Jinlong Yang, Yan Zhang, Michael J. Black, Krikamol Muandet, and Siyu Tang. Grasping Field: Learning Implicit Representations for Human Grasps. In *International Conference on 3D Vision*, 2020. 2, 3, 9

[24] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Khan, and Mubarak Shah. Transformers in Vision: A Survey. In *arXiv Preprint*, 2021. 3

[25] Dong Uk Kim, Kwang In Kim, and Seungryul Baek. End-to-End Detection and Pose Estimation of Two Interacting Hands. In *International Conference on Computer Vision*, 2021. 2, 3, 6

[26] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *International Conference for Learning Representations*, 2015. 10

[27] Dominik Kulon, Riza Alp Güler, Iasonnas Kokkinos, Michael Bronstein, and Stefanos Zafeiriou. Weakly-Supervised Mesh-Convolutional Hand Reconstruction in the Wild. In *Conference on Computer Vision and Pattern Recognition*, 2020. 5

[28] Manoj Kumar, Dirk Weissenborn, and Nal Kalchbrenner. Colorization Transformer. In *arXiv Preprint*, 2021. 3

[29] Taein Kwon, Bugra Tekin, Jan Stühmer, Federica Bogo, and Marc Pollefeys. H2O: Two Hands Manipulating Objects for First Person Interaction Recognition. In *International Conference on Computer Vision*, 2021. 3

[30] Nikolaos Kyriazis and Antonis A. Argyros. Scalable 3D Tracking of Multiple Interacting Objects. In *Conference on Computer Vision and Pattern Recognition*, 2014. 2

[31] Kailin Li, Lixin Yang, Xinyu Zhan, Jun Lv, Wenqiang Xu, Jiefeng Li, and Cewu Lu. ArtiBoost: Boosting articulated 3d hand-object pose estimation via online exploration and

[32] synthesis. In *Conference on Computer Vision and Pattern Recognition*, 2022. 9, 10

[32] Kevin Lin, Lijuan Wang, and Zicheng Liu. End-to-End Human Pose and Mesh Reconstruction with Transformers. In *Conference on Computer Vision and Pattern Recognition*, 2021. 3, 10

[33] Kevin Lin, Lijuan Wang, and Zicheng Liu. Mesh Graphormer. In *International Conference on Computer Vision*, 2021. 3

[34] Shaowei Liu, Hanwen Jiang, Jiarui Xu, Sifei Liu, and Xiaolong Wang. Semi-Supervised 3D Hand-Object Poses Estimation with Interactions in Time. In *Conference on Computer Vision and Pattern Recognition*, 2021. 10

[35] Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. Camera Distance-Aware Top-Down Approach for 3D Multi-Person Pose Estimation from a Single RGB Image. In *International Conference on Computer Vision*, 2019. 9

[36] Gyeongsik Moon and Kyoung Mu Lee. I2L-MeshNet: Image-to-Lixel Prediction Network for Accurate 3D Human Pose and Mesh Estimation from a Single RGB Image. In *European Conference on Computer Vision*, 2020. 1, 10

[37] 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 *European Conference on Computer Vision*, 2020. 1, 2, 3, 5, 6, 7, 8, 9, 11, 14, 15

[38] Franziska Mueller, Micah Davis, Florian Bernard, Oleksandr Sotnychenko, Mickeal Verschoor, Miguel Otaduy, Dan Casas, and Christian Theobalt. Real-Time Pose and Shape Reconstruction of Two Interacting Hands with a Single Depth Camera. *IEEE Transactions on Robotics and Automation*, 38, 2019. 2, 3

[39] Markus Oberweger, Paul Wohlhart, and Vincent Lepetit. Training a Feedback Loop for Hand Pose Estimation. In *International Conference on Computer Vision*, 2015. 1

[40] Iason Oikonomidis, Nikolaos Kyriazis, and Antonis A. Argyros. Full DOF Tracking of a Hand Interacting with an Object by Modeling Occlusions and Physical Constraints. In *International Conference on Computer Vision*, 2011. 2

[41] I. Oikonomidis, Nikolaos Kyriazis, and Antonis A. Argyros. Tracking the Articulated Motion of Two Strongly Interacting Hands. In *Conference on Computer Vision and Pattern Recognition*, 2012. 2

[42] Paschalis Panteleris, Nikolaos Kyriazis, and Antonis A. Argyros. 3D Tracking of Human Hands in Interaction with Unknown Objects. In *British Machine Vision Conference*, 2015. 2

[43] Paschalis Panteleris, Iason Oikonomidis, and Antonis A. Argyros. Using a Single RGB Frame for Real Time 3D Hand Pose Estimation in the Wild. In *IEEE Winter Conference on Applications of Computer Vision*, 2018. 1, 2, 3

[44] JoonKyu Park, Yeonguk Oh, Gyeongsik Moon, Hongsuk Choi, and Kyoung Mu Lee. HandOccNet: Occlusion-Robust 3D Hand Mesh Estimation Network. In *Conference on Computer Vision and Pattern Recognition*, 2022. 9, 10

[45] Kiru Park, Timothy Patten, and Markus Vincze. Pix2Pose: Pixel-Wise Coordinate Regression of Objects for 6D Pose

Estimation. In *International Conference on Computer Vision*, 2019. 5

[46] Georgios Pavlakos, Nikos Kolotouros, and Kostas Daniilidis. TexturePose: Supervising Human Mesh Estimation with Texture Consistency. In *International Conference on Computer Vision*, 2019. 5, 8

[47] Georgios Pavlakos, Xiaowei Zhou, Aaron Chan, Konstantinos G Derpanis, and Kostas Daniilidis. 6-dof object pose from semantic keypoints. In *International Conference on Robotics and Automation (ICRA)*, 2017. 9

[48] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to Estimate 3D Human Pose and Shape from a Single Color Image. In *Conference on Computer Vision and Pattern Recognition*, 2018. 1, 5, 8

[49] Javier Romero, Dimitrios Tzionas, and Michael J. Black. EMbodied Hands: Modeling and Capturing Hands and Bodies Together. *IEEE Transactions on Robotics and Automation*, 36, 2017. 2, 3, 5, 9

[50] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *Conference on Medical Image Computing and Computer Assisted Intervention*, 2015. 4, 9

[51] Breannan Smith, Chenglei Wu, He Wen, Patrick Peluse, Yaser Sheikh, Jessica Hodgins, and Takaaki Shiratori. Constraining Dense Hand Surface Tracking with Elasticity. *IEEE Transactions on Robotics and Automation*, 39, 2020. 2

[52] Adrian Spurr, Umar Iqbal, Pavlo Molchanov, Otmar Hilliges, and Jan Kautz. Weakly Supervised 3D Hand Pose Estimation via Biomechanical Constraints. In *European Conference on Computer Vision*, 2020. 1

[53] Srinath Sridhar, Franziska Mueller, Michael Zollhöfer, Dan Casas, Antti Oulasvirta, and Christian Theobalt. Real-Time Joint Tracking of a Hand Manipulating an Object from RGB-D Input. In *European Conference on Computer Vision*, 2016. 2

[54] Omid Taheri, Nima Ghorbani, Michael J. Black, and Dimitrios Tzionas. GRAB: A Dataset of Whole-Body Human Grasping of Objects. In *European Conference on Computer Vision*, 2020. 3, 5

[55] Jonathan Taylor, Lucas Bordeaux, Thomas Cashman, Bob Corish, Cem Keskin, Toby Sharp, Eduardo Soto, David Sweeney, Julien Valentin, Benjamin Luff, Arran Topalian, Erroll Wood, Sameh Khamis, Pushmeet Kohli, Shahram Izadi, Richard Banks, Andrew Fitzgibbon, and Jamie Shotton. Efficient and Precise Interactive Hand Tracking through Joint, Continuous Optimization of Pose and Correspondences. *IEEE Transactions on Robotics and Automation*, 35, 2016. 2

[56] Jonathan Taylor, Vladimir Tankovich, Danhang Tang, Cem Keskin, David Kim, Philip Davidson, Adarsh Kowdle, and Shahram Izadi. Articulated Distance Fields for Ultra-Fast Tracking of Hands Interacting. *IEEE Transactions on Robotics and Automation*, 36, 2017. 2

[57] Bugra Tekin, Federica Bogo, and Marc Pollefeys. H+O: Unified Egocentric Recognition of 3D Hand-Object Poses and Interactions. In *Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2, 3

[58] Hsiao-Yu Tung, Hsiao-Wei Tung, Ersin Yumer, and Katerina Fragkiadaki. Self-Supervised Learning of Motion Capture. In *Advances in Neural Information Processing Systems*, 2017. 1

[59] Dimitrios Tzionas, Luca Ballan, Abhilash Srikantha, Pablo Aponte, Marc Pollefeys, and Juergen Gall. Capturing Hands in Action Using Discriminative Salient Points and Physics Simulation. *International Journal of Computer Vision*, 118, 2016. 2

[60] Dimitrios Tzionas and Juergen Gall. 3D Object Reconstruction from Hand-Object Interactions. In *International Conference on Computer Vision*, 2015. 2

[61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *Advances in Neural Information Processing Systems*, 2017. 4, 5

[62] Jiayi Wang, Franziska Mueller, Florian Bernard, Suzanne Sorli, Oleksandr Sotnychenko, Neng Qian, Miguel Otaduy, Dan Casas, and Christian Theobalt. RGB2Hands: Real-Time Tracking of 3D Hand Interactions from Monocular RGB Video. *IEEE Transactions on Robotics and Automation*, 39, 2020. 1, 2, 3

[63] Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. Monocular Total Capture: Posing Face, Body, and Hands in the Wild. In *Conference on Computer Vision and Pattern Recognition*, 2019. 2

[64] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. PoseCNN: A Convolutional Neural Network for 6D Object Pose Estimation in Cluttered Scenes. *Science*, 2018. 3, 6, 13

[65] Fuzhi Yang, Huan Yang, Jianlong Fu, Hongtao Lu, and Baining Guo. Learning Texture Transformer Network for Image Super-Resolution. In *Conference on Computer Vision and Pattern Recognition*, 2020. 3

[66] 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 *International Conference on Computer Vision*, 2021. 2, 3

[67] Xiaozheng Zheng, Pengfei Ren, Haifeng Sun, Jingyu Wang, Qi Qi, and Jianxin Liao. Joint-aware regression: Rethinking regression-based method for 3d hand pose estimation. In *British Machine Vision Conference*, 2021. 10

[68] Y. Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the Continuity of Rotation Representations in Neural Networks. In *Conference on Computer Vision and Pattern Recognition*, 2019. 5

[69] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: Deformable Transformers for End-to-End Object Detection. In *International Conference for Learning Representations*, 2021. 3

[70] Christian Zimmermann and Thomas Brox. Learning to Estimate 3D Hand Pose from Single RGB Images. In *International Conference on Computer Vision*, 2017. 1

[71] Christian Zimmermann, Duygu Ceylan, Jimei Yang, Bryan C. Russell, Max Argus, and Thomas Brox. FreiHAND: A Dataset for Markerless Capture of Hand Pose and

Shape from Single RGB Images. In *International Conference on Computer Vision*, 2019. 2, 3