

TOCH: Spatio-Temporal Object-to-Hand Correspondence for Motion Refinement

Keyang Zhou^{1,2}, Bharat Lal Bhatnagar^{1,2},
Jan Eric Lenssen², and Gerard Pons-Moll^{1,2}

¹ University of Tübingen, Germany

{keyang.zhou,gerard.pons-moll}@uni-tuebingen.de

² Max Planck Institute for Informatics, Saarland Informatics Campus, Germany

{bhatnag,jlenssen}@mpi-inf.mpg.de

Abstract. We present TOCH, a method for refining incorrect 3D hand-object interaction sequences using a correspondence based prior learnt directly from data. Existing hand trackers, especially those that rely on very few cameras, often produce visually unrealistic results with hand-object intersection or missing contacts. Although correcting such errors requires reasoning about temporal aspects of interaction, most previous works focus on static grasps and contacts. The core of our method are TOCH fields, a novel spatio-temporal representation for modeling correspondences between hands and objects during interaction. TOCH fields are a point-wise, object-centric representation, which encode the hand position relative to the object. Leveraging this novel representation, we learn a latent manifold of plausible TOCH fields with a temporal denoising auto-encoder. Experiments demonstrate that TOCH outperforms state-of-the-art 3D hand-object interaction models, which are limited to static grasps and contacts. More importantly, our method produces smooth interactions even before and after contact. Using a single trained TOCH model, we quantitatively and qualitatively demonstrate its usefulness for correcting erroneous sequences from off-the-shelf RGB/RGB-D hand-object reconstruction methods and transferring grasps across objects. Our code and model are available at [1].

Keywords: hand-object interaction, motion refinement, hand prior

1 Introduction

Tracking hands that are in interaction with objects is an important part of many applications in Virtual and Augmented Reality, such as modeling digital humans capable of manipulation tasks [59, 23, 77, 6, 70]. Although there exists a vast amount of literature about tracking hands in isolation, much less work has focused on joint tracking of objects and hands. The high degrees of freedom in possible hand configurations, frequent occlusions, noisy or incomplete observations (*e.g.* lack of depth channel in RGB images) make the problem heavily ill-posed. We argue that tracking interacting hands requires a powerful prior

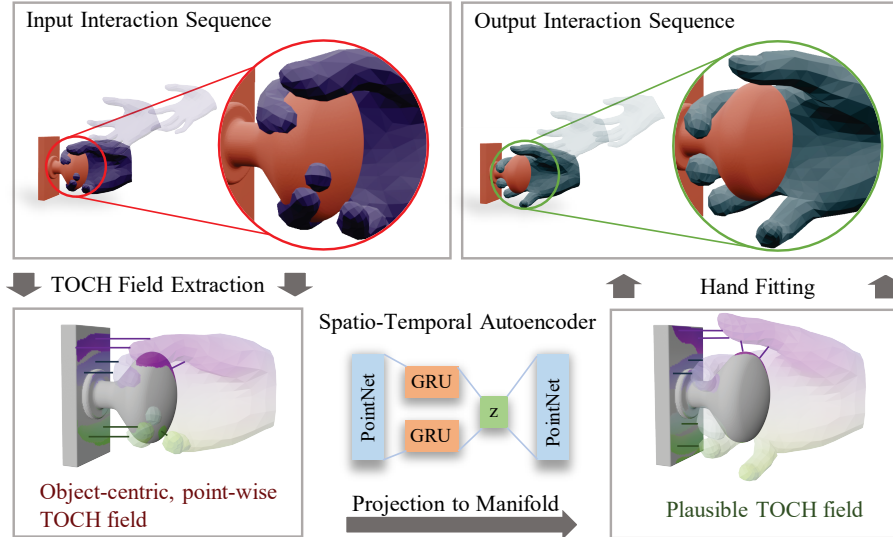


Fig. 1: We propose TOCH, a model for correcting erroneous hand-object interaction sequences. TOCH takes as input a tracking sequence produced by any existing tracker. We extract TOCH fields, a novel object-centric correspondence representation, from a noisy hand-object mesh sequence. The extracted noisy TOCH fields are fed into an auto-encoder, which projects it onto a learned hand motion manifold. Lastly, we obtain the corrected tracking sequence by fitting hands to the reconstructed TOCH fields. TOCH is applicable to interaction sequences even before and after contact happens.

learned from a set of clean interaction sequences, which is the core principle of our method.

Beyond the aforementioned challenges, subtle errors in hand estimation have a huge impact on perceived realism. For example, if the 3D object is floating in the air, is grasped in a non-physically plausible way, or hand and object intersect, the perceived quality will be poor. Unfortunately, such artifacts are common in pure hand-tracking methods. Researchers have used different heuristics to improve plausibility, such as inter-penetration constraints [28] and smoothness priors [26]. A recent line of work predicts likely static hand poses and grasps for a given object [35, 24] but those methods can not directly be used as a prior to fix common capturing and tracking errors. Although there exists work to refine hand-object interactions [60, 22], it is only concerned with static grasps.

In this work, we propose TOCH, a data-driven method for refining noisy 3D hand-object interaction sequences. In contrast to previous work in interaction modeling, TOCH not only considers static interactions but can also be applied to sequences without introducing snapping artifacts. The whole approach is outlined in Figure 1. Our key insight is that estimating point-wise, spatio-temporal

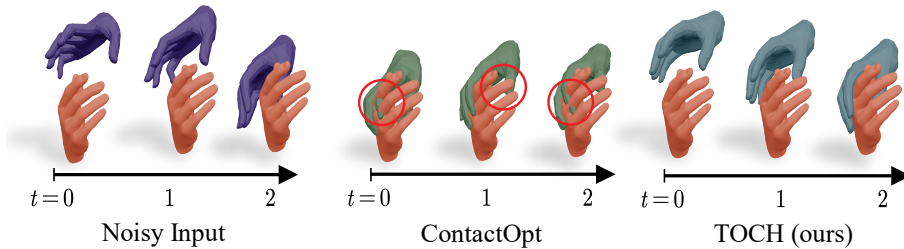


Fig. 2: Example refinement of an interaction sequence. Left: a noisy sequence of a hand approaching and grasping another static hand. Middle: ContactOpt [22] always snaps the hand into grasping posture regardless of its position relative to the object, as it is not designed for sequences. Right: TOCH preserves the relative hand-object arrangement during interaction while refining the final grasp.

object-hand correspondences are crucial for realism, and sufficient to constrain the high-dimensional hand poses. Thus, the point-wise correspondences between object and hand are encoded in a novel spatio-temporal representation called *TOCH field*, which takes the object geometry and the configuration of the hand with respect to the object into account. We then learn the manifold of plausible TOCH fields from the recently released MoCap dataset of hand-object interactions [60] using an auto-encoder and apply it to correcting noisy observations. In contrast to conventional binary contacts [11, 22, 67], TOCH fields also encode the position of hand parts that are not directly in contact with the object, making TOCH applicable to whole interaction sequences, see Figure 2. TOCH has further useful properties for practical application:

- TOCH can effectively project implausible hand motions to the learned object-centric hand motion manifold and produces visually correct interaction sequences that outperform previous static approaches.
- TOCH does not depend on specific sensor data (RGB image, depth map *etc.*) and can be used with any tracker.
- TOCH can be used to transfer grasp sequences across objects sharing similar geometry, even though it is not designed for this task.

2 Related Work

2.1 Hand and Object Reconstruction

Hand Reconstruction and Tracking. Reconstructing 3D hand surfaces from RGB or depth observations is a well-studied problem [31]. Existing work can generally be classified into two paradigms: discriminative approaches [20, 14, 80, 44, 9] directly estimate hand shape and pose parameters from the observation, while generative approaches [58, 61, 63] iteratively optimize a parametric hand

model so that its projection matches the observation. Recently, more challenging settings such as reconstructing two interacting hands [72, 47, 56] are also explored. These works ignore the presence of objects and are hence less reliable in interaction-intensive scenarios.

Joint Hand and Object Reconstruction. Jointly reconstructing hand and object in interaction [49, 4, 66, 57, 73, 74, 51] has received much attention. Owing to the increasing amount of hand-object interaction datasets with annotations [29, 25, 11, 19, 82, 39], deep neural networks are often used to estimate an initial hypothesis for hand and object poses, which are then jointly optimized to meet certain interaction constraints [29, 26, 15, 27, 13]. Most works in this direction improve contact realism by encouraging a small hand-to-object distance and penalizing inter-penetrating vertices [33, 4]. However, these simple approaches often yield implausible interaction and do not take the whole motion sequence into account. In contrast, our method alleviates both shortcomings through a object-centric, temporal representation that also considers frames in which hand and object are not in direct contact.

2.2 Hand Contact and Grasp

Grasp Synthesis. Synthesizing novel hand grasp given an object has been widely studied in robotics [55]. Traditional approaches either optimize for force-closure condition [17] or sample and rank grasp candidates based on learned features [8]. There are also hybrid approaches that combine the merits of both [45, 40]. Recently, a number of neural network-based models have been proposed for this task [28, 60, 16, 81, 32]. In particular, [35, 34] represent the hand-object proximity as an implicit function. We took a similar approach and represent the hand relative to the object by signed distance values distributed on the object.

Object Manipulation Synthesis. In comparison with static grasp synthesis, generating dexterous manipulation of objects is a more difficult problem since it additionally requires dynamic hand and object interaction to be modeled. This task is usually approached by optimizing hand poses to satisfy a range of contact force constraints [42, 69, 46, 79]. Hand motions generated by these works are physically plausible but lack natural variations. Zhang *et al.* [75] utilized various hand-object spatial representations to learn object manipulation from data. An IK solver is used to avoid inter-penetration. We took a different approach and solely use an object-centric spatio-temporal representation, which is shown to be less prone to interaction artifacts.

Contact Refinement. Recently, some works focus on refining hand and object contact [60, 68, 22]. Both [68] and [22] propose to first estimate the potential contact region on the object and then fit the hand to match the predicted contact. However, limited by the proposed contact representation, they can only model hand and object *in stable grasp*. While we share a similar goal, our work can also deal with the case where the hand is *close to* but not in contact with the object, as a result of our novel hand-object correspondence representation. Hence our method can be used to refine a tracking sequence.

2.3 Pose and Motion Prior

It has been observed that most human activities lie on low-dimensional manifolds [18, 65]. Therefore natural motion patterns can be found by applying learned data priors. A pose or motion prior can facilitate a range of tasks including pose estimation from images or videos [7, 3, 43, 71], motion interpolation [41], motion capture [76, 64], and motion synthesis [30, 2, 12]. Early attempts in capturing pose and motion priors mostly use simple statistical models such as PCA [50], Gaussian Mixture Models [7] or Gaussian Process Dynamical Models [65]. With the advent of deep generative models [36, 21], recent works rely on auto-encoders [52, 37] and adversarial discriminators [78, 38] to more faithfully capture the motion distribution.

Compared to body motion prior, there is less work devoted to hand motion priors. Ng *et al.* [48] learned a prior of conversational hand gestures conditioned on body motion. Our work bears the most similarity to [24], where an object-dependent hand pose prior was learned to foster tracking. Hamer *et al.* [24] proposed to map hand parts into local object coordinates and learn the object-dependent distribution with a Parzen density estimator. The prior is learned on a few objects and subsequently transferred to objects from the same class by geometric warping. Hence it cannot truly capture the complex correlation between hand gesture and object geometry.

3 Method

In this section, we describe our method for refining hand pose sequences during interaction with an object. We begin by introducing the problem setting and outlining our approach. Let $\mathbf{H} = (\mathbf{H}^i)_{1 \leq i \leq T}$ with $\mathbf{H}^i \in \mathbb{R}^{K \times 3}$ denote a sequence of vertices that describe hand meshes over the course of an interaction over T frames. We only deal with sequences containing a single hand and a single rigid object mesh, whose vertices we denote as $\mathbf{O} \in \mathbb{R}^{L \times 3}$. We assume the object shape to be known. Since we care about hand motion relative to the object, we express the hands in local object space, and the object coordinates remain fixed over the sequence. The per-frame hand vertices \mathbf{H}^i in object space are produced by a parametric hand model MANO [54] using linear blend skinning:

$$\mathbf{H}^i = \text{LBS}(\mathbf{Y}; \boldsymbol{\beta}, \boldsymbol{\theta}^i) + \mathbf{t}_H^i. \quad (1)$$

where the parameters $\{\boldsymbol{\beta}^i, \boldsymbol{\theta}^i, \mathbf{t}^i\}$ denote shape, pose and translation w.r.t. template hand mesh \mathbf{Y} respectively.

Observing the hand-object motion through RGB or depth sensors, a hand tracker yields an estimated hand motion sequence $\tilde{\mathbf{H}} = (\tilde{\mathbf{H}}^i)_{1 \leq i \leq T}$. The goal of our method is to improve the perceptual realism of this potentially noisy estimate using prior information learned from training data.

Concept. We observe that during hand-object interactions, the hand motion is heavily constrained by the object shape. Therefore, noisy hand-object interaction

is a deviation from a low-dimensional manifold of realistic hand motions, conditioned on the object. We formulate our goal as learning a mapping to maximize the posterior $p(\mathbf{H}|\tilde{\mathbf{H}}, \mathbf{O})$ of the real motion \mathbf{H} given the noisy observation $\tilde{\mathbf{H}}$ and the object with which the hand interacts. This amounts to finding an appropriate sequence of MANO parameters, which is done in three steps (see Figure 1): **1)** The initial estimate of a hand motion sequence is encoded with a TOCH field, our object-centric, point-wise correspondence representation (Section 3.1). **2)** The TOCH fields are projected to a learned low-dimensional manifold using a temporal denoising auto-encoder (Section 3.2). **3)** A sequence of corrected hand meshes is obtained from the processed TOCH fields (Section 3.3).

3.1 TOCH Fields

Naively training an auto-encoder on hand meshes is problematic, because the model could ignore the conditioning object and learn a plain hand motion prior (Sec. 4.5). Moreover, if we include the object into the formulation, the model would need to learn manifolds for all joint rigid transformation of hand and object, which leads to high problem complexity [35]. Thus, we represent the hand as a TOCH field \mathbf{F} , which is a spatio-temporal object-centric representation that makes our approach invariant to joint hand and object rotation and translation.

TOCH Field Representation. For an initial estimation $\tilde{\mathbf{H}}$ of the hand mesh and the given object mesh \mathbf{O} , we define the TOCH field as a collection of point-wise vectors on a set $\{\mathbf{o}_i\}_{i=1}^N$ of N points, sampled from the object surface:

$$\mathbf{F}(\tilde{\mathbf{H}}, \mathbf{O}) = \{(c_i, d_i, \mathbf{y}_i)\}_{i=1}^N, \quad (2)$$

where $c_i \in \{0, 1\}$ is a binary flag indicating whether the i -th sampled object point has a corresponding point on the hand surface, $d_i \in \mathbb{R}$ is the signed distance between the object point and its corresponding hand point, and $\mathbf{y}_i \in \mathbb{R}^3$ are the coordinates of the corresponding hand point on the un-posed canonical MANO template mesh. Note that \mathbf{y}_i is a 3D location on the hand surface embedded in \mathbb{R}^3 , encoding the correspondence similar to [5, 62].

Finding correspondences. As we model whole interaction sequences, including frames in which the hand and the object are not in contact, we cannot simply define the correspondences as points that lie within a certain distance to each other. Instead, we generalize the notion of contact by diffusing the object mesh into \mathbb{R}^3 . We cast rays from the object surface along its normal directions, as outlined in Figure 1. The object normal vectors are obtained from the given object mesh. The correspondence of an object point is obtained as the first intersection with the hand mesh. If there is no intersection, or the first intersection is not the hand, this object point has no correspondence. If the object point is inside the hand, which might happen in case of noisy observations, we search for correspondences along the negative normal direction. The detailed procedure for determining correspondences is listed in Algorithm 1.

Algorithm 1: Finding object-hand correspondences

Input: Hand mesh \mathbf{H} , object mesh \mathbf{O} , uniformly sampled object points and normals $\{\mathbf{o}_i, \mathbf{n}_i\}_{i=1}^N$

Output: Binary correspondence indicators $\{c_i\}_{i=1}^N$

for $i = 1$ **to** N **do**

$c_i \leftarrow 0;$
if \mathbf{o}_i inside \mathbf{H} $s \leftarrow -1$ else $s \leftarrow 1$;
$\mathbf{r}_1 \leftarrow \text{ray}(\mathbf{o}_i, s\mathbf{n}_i);$
$\mathbf{p}_1 \leftarrow \text{ray_mesh_intersection}(\mathbf{r}_1, \mathbf{H});$
if $\mathbf{p}_1 \neq \emptyset$
$\mathbf{r}_2 \leftarrow \text{ray}(\mathbf{o}_i + \epsilon s\mathbf{n}_i, s\mathbf{n}_i);$
$\mathbf{p}_2 \leftarrow \text{ray_mesh_intersection}(\mathbf{r}_2, \mathbf{O});$
if $\mathbf{p}_2 = \emptyset$ or $\ \mathbf{o}_i - \mathbf{p}_1\ < \ \mathbf{o}_i - \mathbf{p}_2\ $
$c_i \leftarrow 1;$

Representation properties. The described TOCH field representation has the following advantages. **1)** It is naturally invariant to joint rotation and translation of object and hand, which reduces required model complexity. **2)** By specifying the distance between corresponding points, TOCH fields enable a subsequent auto-encoder to reason about point-wise proximity of hand and object. This helps to correct various artifacts, *e.g.* inter-penetration can be simply detected by finding object vertices with a negative correspondence distance. **3)** From surface normal directions of object points and the corresponding distances, a TOCH field can be seen as an encoding of the partial hand point cloud from the perspective of the object surface. We can explicitly derive that point cloud from the TOCH field and use it to infer hand pose and shape by fitting the hand model to the point cloud (c.f. Section 3.3).

3.2 Temporal Denoising Auto-encoder

To project a noisy TOCH field to the correct manifold, we use a temporal denoising auto-encoder, consisting of an encoder $g_{\text{enc}} : (\tilde{\mathbf{F}}_i)_{1 \leq i \leq T} \mapsto (\mathbf{z}_i)_{1 \leq i \leq T}$, which maps a sequence of noisy TOCH fields (concatenated with the coordinates and normals of each object point) to latent representation, and a decoder $g_{\text{dec}} : (\mathbf{z}_i)_{1 \leq i \leq T} \mapsto (\hat{\mathbf{F}}_i)_{1 \leq i \leq T}$, which computes the corrected TOCH fields $\hat{\mathbf{F}}$ from the latent codes. As TOCH fields consist of feature vectors attached to points, we use a PointNet-like [53] architecture. The point features in each frame are first processed by consecutive PointNet blocks to extract frame-wise features. These features are then fed into a bidirectional GRU layer to capture temporal motion patterns. The decoder network again concatenates the encoded frame-wise features with coordinates and normals of the object points and produces denoised TOCH fields $(\hat{\mathbf{F}}_i)_{1 \leq i \leq T}$. The network is trained by minimizing

$$\mathcal{L}(\hat{\mathbf{F}}, \mathbf{F}) = \sum_{i=1}^T \sum_{j=1}^N c_j^i \left(\|\hat{\mathbf{y}}_j^i - \mathbf{y}_j^i\|_2^2 + w_{ij} (\hat{d}_j^i - d_j^i)^2 \right) - \text{BCE}(\hat{c}_j^i, c_j^i), \quad (3)$$

where \mathbf{F} denotes the groundtruth TOCH fields and $\text{BCE}(\hat{c}_j^i, c_j^i)$ is the binary cross entropy between output and target correspondence indicators. Note that we only compute the first two parts of the loss on TOCH field elements with $c_j^i = 1$, i.e. object points that have a corresponding hand point. We use a weighted loss on the distances \hat{d}_j^i . The weights are defined as

$$w_{ij} = \frac{\exp(-\|\hat{d}_j^i\|)}{\sum_{k=1}^{N_i} \exp(-\|d_k^i\|)} N_i, \quad (4)$$

where $N_i = \sum_{j=1}^N c_j^i$. This weighting scheme encourages the network to focus on regions of close interaction, where a slight error could have huge impact on contact realism. Multiplying by the sum of correspondence ensures equal influence of all points in the sequence instead of equal influence of all frames.

3.3 Hand Motion Reconstruction

After projecting the noisy TOCH fields of input tracking sequence to the manifold learned by the auto-encoder, we need to recover the hand motion from the processed TOCH fields. The TOCH field is not fully differentiable w.r.t. the hand parameters, as changing correspondences would involve discontinuous function steps. Thus, we cannot directly optimize the hand pose parameters to produce the target TOCH field. Instead, we decompose the optimization into two steps. We first use the denoised TOCH fields to locate hand points corresponding to the object points. We then optimize the MANO model to find hands that best fit these points, which is a differentiable formulation.

Formally, given denoised TOCH fields $\mathbf{F}^i(\mathbf{H}, \mathbf{O}) = \{(c_j^i, d_j^i, \mathbf{y}_j^i)\}_{j=1}^N$ for frames $i \in \{1, \dots, T\}$ on object points $\{\mathbf{o}_j\}_{j=1}^N$, we first produce the partial point clouds $\hat{\mathbf{Y}}^i$ of the hand as seen from the object’s perspective:

$$\hat{\mathbf{y}}_j^i = \mathbf{o}_j + d_j^i \mathbf{n}_j^i. \quad (5)$$

Then, we fit MANO to those partial point clouds by minimizing:

$$\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\theta}, \mathbf{t}_H) = \sum_{i=1}^T \mathcal{L}_{\text{corr}}(\boldsymbol{\beta}, \boldsymbol{\theta}^i, \mathbf{t}_H) + \mathcal{L}_{\text{reg}}(\boldsymbol{\beta}, \boldsymbol{\theta}). \quad (6)$$

The first term of Equation 6 is the hand-object correspondence loss

$$\mathcal{L}_{\text{corr}}(\boldsymbol{\beta}, \boldsymbol{\theta}^i, \mathbf{t}_H) = \sum_{j=1}^N c_j^i \left\| \hat{\mathbf{y}}_j^i - (\text{LBS}(\text{Proj}_{\mathbf{Y}}(\mathbf{y}_j^i); \boldsymbol{\beta}, \boldsymbol{\theta}^i) + \mathbf{t}_H) \right\|^2, \quad (7)$$

where LBS is the linear blend skinning function in Equation 1 and $\text{Proj}_{\mathbf{Y}}(\cdot)$ projects a point to the nearest point on the template hand surface. This loss term ensures that the deformed template hand point corresponding to \mathbf{o}_i is at a predetermined position derived from the TOCH field.

The last term of (6) regularizes shape and pose parameters of MANO,

$$\mathcal{L}_{\text{reg}}(\boldsymbol{\beta}, \boldsymbol{\theta}) = w_1 \|\boldsymbol{\beta}\|^2 + w_2 \sum_{i=1}^T \|\boldsymbol{\theta}^i\|^2 + w_3 \sum_{i=1}^{T-1} \|\boldsymbol{\theta}^{i+1} - \boldsymbol{\theta}^i\|^2 + w_4 \sum_{i=2}^{T-1} \sum_{k=1}^J \|\ddot{\mathbf{p}}_k^i\| \quad (8)$$

where $\ddot{\mathbf{p}}_k^i$ is the acceleration of hand joint k in frame i . Besides regularizing the norm of MANO parameters, we additionally enforce temporal smoothness of hand poses. This is necessary because (7) only constrains those parts of a hand with object correspondences. Per-frame fitting of TOCH fields leads to multiple plausible solutions, which can only be disambiguated by considering neighbouring frames. Since (6) is highly nonconvex, we optimize it in two stages. In the first stage, we freeze shape and pose, and only optimize hand orientation and translation. We then jointly optimize all the variables in the second stage.

4 Experiments

In this section, we evaluate the presented method on synthetic and real datasets of hand/object interaction. Our goal is to verify that TOCH produces *realistic interaction sequences* (Section 4.3), *outperforms previous static approaches* in several metrics (Section 4.4), and derives a *meaningful representation* for hand object interaction (Section 4.5). Before presenting the results, we introduce the used datasets in Section 4.1 and the evaluated metrics in Section 4.2.

4.1 Datasets

GRAB. We train TOCH on GRAB [60], a MoCap dataset for whole-body grasping of objects. GRAB contains interaction sequences with 51 objects from [10]. We pre-select 10 objects for validation and testing, and train with the rest sequences. Since we are only interested in frames where interaction is about to take place, we filter out frames where the hand wrist is more than 15 cm away from the object. Due to symmetry of the two hands, we anchor correspondences to the right hand and flip left hands to increase the amount of training data.

HO-3D. HO-3D is a dataset of hand-object video sequences captured by RGB-D cameras. It provides frame-wise annotations for 3D hand poses and 6D object poses, which are obtained from a novel joint optimization procedure. To ensure fair comparison with baselines which are not designed for sequences without contact, we compare on a selected subset of static frames with hand-object contact.

4.2 Metrics

Mean Per-Joint Position Error (MPJPE). We report the average Euclidean distance between refined and groundtruth 3D hand joints. Since pose annotation quality varies across datasets, this metric should be jointly assessed with other perceptual metrics.

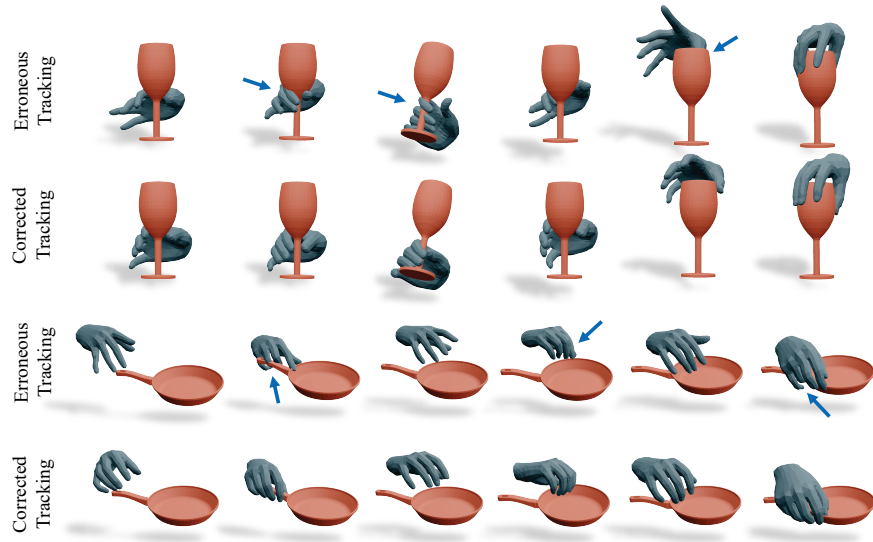


Fig. 3: Qualitative results on two synthetic hand-object interaction sequences that suffer from inter-penetration and non-smooth hand motion. The results after TOCH refinement show correct contact and are much more visually plausible. Note that TOCH only applies minor changes in hand poses but the perceived realism is largely enhanced. Check the supplemental video for animated results.

Mean Per-Vertex Position Error (MPVPE). This metric represents the average Euclidean distance between refined and groundtruth 3D meshes. It assesses the reconstruction accuracy of both hand shape and pose.

Solid Intersection Volume (IV). We measure hand-object inter-penetration by voxelizing hand and object meshes and reporting the volume of voxels occupied by both. Solely considering this metric can be misleading since it does not account for the case where the object is not in effective contact with the hand.

Contact IoU (C-IoU). This metric assesses the Intersection-over-Union between the groundtruth contact map and the predicted contact map. The contact map is obtained from the binary hand-object correspondence by thresholding the correspondence distance within ± 2 mm. We only report this metric on GRAB since the groundtruth annotations in HO-3D are not accurate enough [22].

4.3 Refining Synthetic Tracking Error

In order to use TOCH in real settings, it would be ideal to train the model on the predictions of existing hand trackers. However, this requires large amount of images/depth sequences paired with accurate hand and object annotations, which is currently not available. Moreover, targeting a specific tracker might lead to overfitting to tracker-specific errors, which is undesirable for generalization.

GRAB-Type → Noise →	GRAB-T (0.01)	GRAB-T (0.02)	GRAB-R (0.3)	GRAB-R (0.5)	GRAB-B (0.01 & 0.3)
MPJPE (mm) ↓	16.0 → 9.93	31.9 → 12.3	4.58 → 9.58	7.53 → 9.12	17.3 → 10.3
MPVPE (mm) ↓	16.0 → 11.8	31.9 → 13.9	6.30 → 11.5	10.3 → 11.0	18.3 → 12.1
IV (cm ³) ↓	2.48 → 1.79	2.40 → 2.50	1.88 → 1.52	1.78 → 1.35	2.20 → 1.78
C-IoU (%) ↑	3.56 → 29.2	2.15 → 16.7	11.4 → 26.6	5.06 → 24.4	1.76 → 26.6

Table 1: We quantitatively evaluate TOCH on multiple perturbed GRAB test sets with different types and magnitude of noise. The numbers inside the parentheses indicate standard deviation of the sampled Gaussian noise. Although pose accuracy is not always improved, TOCH significantly boosts interaction realism for all noise levels, which is demonstrated by the increase in contact IoU and reduction in hand-object inter-penetration.

Method	HO-3D		
	MPJPE (mm) ↓	MPVPE (mm) ↓	IV (cm ³) ↓
Hasson <i>et al.</i>	11.4	11.4	9.26
RefineNet	11.6	11.5	8.11
ContactOpt	9.47	9.45	5.71
TOCH (ours)	9.32	9.28	4.66

Table 2: Quantitative evaluation on HO-3D compared to Hasson *et al.* [26], RefineNet [60] and ContactOpt [22]. We follow the evaluation protocol of HO-3D and report hand joint and mesh errors after Procrustes alignment. TOCH outperforms all the baselines in terms of pose error and interaction quality.

We observe that hand errors can be decomposed into inaccurate global translation and inaccurate joint rotations, and the inaccuracies produced by most state-of-the-art trackers are small. Therefore, we propose to synthesize tracking errors by manually perturbing the groundtruth hand poses of the GRAB dataset. To this end, we apply three different types of perturbation to GRAB: translation-dominant perturbation (abbreviated GRAB-T in the table) applies an additive noise to hand translation \mathbf{t}_H only, pose-dominant perturbation (abbreviated GRAB-R) applies an additive noise to hand pose θ only, and balanced perturbation (abbreviated GRAB-B) uses a combination of both. We only train on the last type while evaluate on all three. The quantitative results are shown in Table 1 and qualitative results are presented in Figure 3.

We can make the following observations. First, TOCH is most effective for correcting translation-dominant perturbations of the hand. For pose-dominant perturbations where the vertex and joint errors are already very small, the resulting hands after TOCH refinement exhibit larger errors. This is because TOCH aims to improve interaction quality of a tracking sequence, which can hardly be reflected by distance based metrics such as MPJPE and MPVPE. We argue

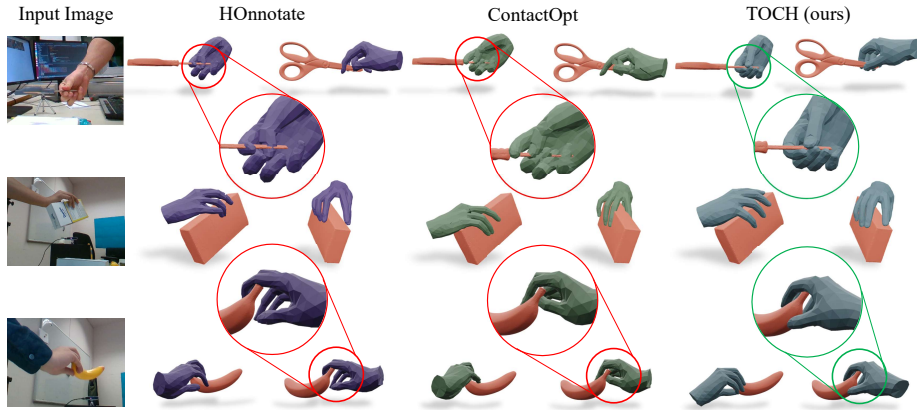


Fig. 4: Qualitative comparison with HOnnotate and ContactOpt. Each sample reconstruction is visualized in two views, the image-aligned view and a side view. We can clearly see hand-object inter-penetrations for HOnnotate and ContactOpt, while our reconstructions are more visually realistic.

that more important metrics for interaction are intersection volume and contact IoU. As an example, the perturbation of GRAB-R (0.3) only induces a tiny joint position error of 4.6 mm, while it results in a significant 88.6% drop in contact IoU. This validates our observation that any slight change in pose has a notable impact on physical plausibility of interaction. TOCH effectively reduces hand-object intersection as well as boosts the contact IoU even when the noise of testing data is higher than that of training data.

4.4 Refining RGB(D)-based Hand Estimators

To evaluate how well TOCH generalizes to real tracking errors, we test TOCH on state-of-the-art models for joint hand-object estimation from image or depth sequences. We first report comparisons with the RGB-based hand pose estimator [26], and two grasp refinement methods RefineNet [60] and ContactOpt [22] in Table 2. Hasson *et al.* [26] predict hand meshes from images, while RefineNet and ContactOpt have no knowledge about visual observations and directly refine hands based on 3D inputs. Groundtruth object meshes are assumed to be given for all the methods. TOCH achieves the best score for all three metrics on HO-3D. In particular, it reduces the mesh intersection volume, indicating an improved interaction quality. We additionally evaluate TOCH on HOnnotate [25], a state-of-the-art RGB-D tracker which annotates the groundtruth for HO-3D. Figure 4 shows some of its failure cases and how they are corrected by TOCH.

4.5 Analysis and Ablation Studies

Grasp transfer. In order to demonstrate the wide-applicability of our learned features, we utilize the pre-trained TOCH auto-encoder for grasp transfer al-

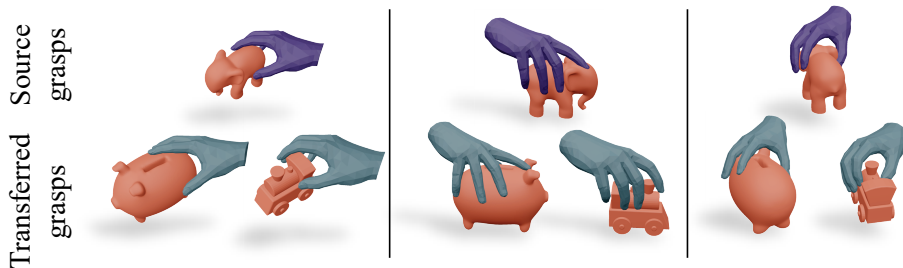


Fig. 5: Transferring grasping poses across objects of different geometry. The top row shows three different source grasps which are subsequently transferred to two target objects in the bottom row. The hand poses are adjusted according to shape of target objects while preserving overall contact.

Method	MPJPE (mm) ↓	IV (cm ³) ↓	C-IoU (%) ↑
Hand-centric baseline	11.2	2.03	18.9
TOCH (w/o corr.)	12.2	2.10	18.6
TOCH (w/o GRU)	10.8	1.87	20.4
TOCH (same obj.)	11.7	1.95	23.1
TOCH (full model)	10.3	1.78	26.6

Table 3: Comparison with various baselines on GRAB-B (0.01 & 0.3). We show that TOCH achieves the lowest hand joint error and intersection volume while recovers the highest percentage of contact regions among all the baselines.

though it was not trained for this task. The goal is to transfer grasping sequences from one object to another object while maintaining plausible contacts. Specifically, given a source hand motion sequence and a source object, we extract the TOCH fields and encode them with our learned encoder network. We then simply decode using the target object – we perform a point-wise concatenation of the latent vectors with point clouds of the target object, and reconstruct TOCH fields with the decoder. This way we can transfer the TOCH fields from the source object to the target object. Qualitative examples are shown in Figure 5.

Object-centric representation. To show the importance of the object-centric representation, we train a baseline model which directly takes noisy hand joint sequences $\{\tilde{\mathbf{j}}^i\}_{i=1}^T$ as input and naively condition it on the object motion sequence $\{\mathbf{O}^i\}_{i=1}^T$. See Table 3 for a quantitative comparison with TOCH. We can observe that although the hand-centric baseline makes small errors in joint positions, the resulting motion is less physically plausible, as reflected by its higher interpenetration and lower contact IoU.

Semantic correspondence. We argue that explicitly reasoning about dense correspondence plays a key role in modeling hand-object interaction. To show this, we train another baseline model in the same manner as in Section 3, except

that we adopt a simpler representation $F(\mathbf{H}, \mathbf{O}) = \{(c_i, d_i)\}_{i=1}^N$, where we keep the binary indicator and signed distance without specifying which hand point is in correspondence. The loss term (7) accordingly changes from mean squared error to Chamfer distance. We can see from Table 3 that this baseline model gives the worst results in all three metrics.

Train and test on the same objects. We test the scenario where objects in test sequences are also seen at training time. We split the dataset based on the action intent label instead of by objects. Specifically, we train on sequences labelled as 'use', 'pass' and 'lift', and evaluate on the remaining. Results from Table. 3 show that generalization to different objects works slightly better than generalizing to different actions. Note that the worse results are also partly attributed to the smaller training set under this new split.

Temporal modeling. We verify the effect of temporal modeling by replacing the GRU layer with global feature aggregation. We concatenate the global average latent code with per-frame latent codes and feed the concatenated feature of each frame to a fully connected layer. As seen in Table. 3, temporal modeling with GRU largely improves interaction quality in terms of recovered contact.

Complexity and running time. The main overhead incurred by TOCH field is in computing ray-triangle intersections, the complexity of which depends on the object geometry and the specific hand-object configuration. As an illustration, it takes around 2s per frame to compute the TOCH field on 2000 sampled object points for an object mesh with 48k vertices and 96k triangles on Intel Xeon CPU. In hand-fitting stage, TOCH is significantly faster than ContactOpt since the hand-object distance can be minimized with mean squared error loss once correspondences are known. Fitting TOCH to a sequence runs at approximately 1 fps on average while it takes ContactOpt over a minute to fit a single frame.

5 Conclusion

In this paper, we introduced TOCH, a spatio-temporal model of hand-object interactions. Our method encodes the hand with TOCH fields, an effective novel object-centric correspondence representation which captures the spatio-temporal configurations of hand and object even before and after contact occurs. TOCH reasons about hand-object configurations beyond plain contacts, and is naturally invariant to rotation and translation. Experiments demonstrate that TOCH outperforms previous methods on the task of 3D hand-object refinement. In future work, we plan to extend TOCH to model more general human-scene interactions.

Acknowledgements This work is supported by the German Federal Ministry of Education and Research (BMBF): Tübingen AI Center, FKZ: 01IS18039A. This work is funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - 409792180 (Emmy Noether Programme, project: Real Virtual Humans). Gerard Pons-Moll is a member of the Machine Learning Cluster of Excellence, EXC number 2064/1 – Project number 390727645. The project was made possible by funding from the Carl Zeiss Foundation.

References

1. <https://virtualhumans.mpi-inf.mpg.de/toch/>
2. Aliakbarian, S., Saleh, F.S., Salzmann, M., Petersson, L., Gould, S.: A stochastic conditioning scheme for diverse human motion prediction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5223–5232 (2020)
3. Arnab, A., Doersch, C., Zisserman, A.: Exploiting temporal context for 3d human pose estimation in the wild. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3395–3404 (2019)
4. Ballan, L., Taneja, A., Gall, J., Van Gool, L., Pollefeys, M.: Motion capture of hands in action using discriminative salient points. In: European Conference on Computer Vision. pp. 640–653. Springer (2012)
5. Bhatnagar, B.L., Sminchisescu, C., Theobalt, C., Pons-Moll, G.: Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. *Advances in Neural Information Processing Systems* **33** (2020)
6. Bhatnagar, B.L., Xie, X., Petrov, I., Sminchisescu, C., Theobalt, C., Pons-Moll, G.: Behave: Dataset and method for tracking human object interactions. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2022)
7. Bogo, F., Kanazawa, A., Lassner, C., Gehler, P., Romero, J., Black, M.J.: Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In: European conference on computer vision. pp. 561–578. Springer (2016)
8. Bohg, J., Morales, A., Asfour, T., Kragic, D.: Data-driven grasp synthesis—a survey. *IEEE Transactions on Robotics* **30**(2), 289–309 (2013)
9. Boukhayma, A., Bem, R.d., Torr, P.H.: 3d hand shape and pose from images in the wild. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10843–10852 (2019)
10. Brahmbhatt, S., Ham, C., Kemp, C.C., Hays, J.: Contactdb: Analyzing and predicting grasp contact via thermal imaging. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8709–8719 (2019)
11. Brahmbhatt, S., Tang, C., Twigg, C.D., Kemp, C.C., Hays, J.: Contactpose: A dataset of grasps with object contact and hand pose. In: *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII* 16. pp. 361–378. Springer (2020)
12. Cai, Y., Wang, Y., Zhu, Y., Cham, T.J., Cai, J., Yuan, J., Liu, J., Zheng, C., Yan, S., Ding, H., et al.: A unified 3d human motion synthesis model via conditional variational auto-encoder. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 11645–11655 (2021)
13. Cao, Z., Radosavovic, I., Kanazawa, A., Malik, J.: Reconstructing hand-object interactions in the wild. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 12417–12426 (2021)
14. Chen, L., Lin, S.Y., Xie, Y., Lin, Y.Y., Xie, X.: Mvbm: A large-scale multi-view hand mesh benchmark for accurate 3d hand pose estimation. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 836–845 (2021)
15. Chen, Y., Tu, Z., Kang, D., Chen, R., Bao, L., Zhang, Z., Yuan, J.: Joint hand-object 3d reconstruction from a single image with cross-branch feature fusion. *IEEE Transactions on Image Processing* **30**, 4008–4021 (2021)

16. Corona, E., Pumarola, A., Alenya, G., Moreno-Noguer, F., Rogez, G.: Ganhand: Predicting human grasp affordances in multi-object scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5031–5041 (2020)
17. El-Khoury, S., Sahbani, A., Bidaud, P.: 3d objects grasps synthesis: A survey. In: 13th World Congress in Mechanism and Machine Science. pp. 573–583 (2011)
18. Elgammal, A., Lee, C.S.: The role of manifold learning in human motion analysis. In: Human Motion, pp. 25–56. Springer (2008)
19. Garcia-Hernando, G., Yuan, S., Baek, S., Kim, T.K.: First-person hand action benchmark with rgb-d videos and 3d hand pose annotations. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 409–419 (2018)
20. Ge, L., Ren, Z., Li, Y., Xue, Z., Wang, Y., Cai, J., Yuan, J.: 3d hand shape and pose estimation from a single rgb image. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10833–10842 (2019)
21. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. *Advances in neural information processing systems* **27** (2014)
22. Grady, P., Tang, C., Twigg, C.D., Vo, M., Brahmbhatt, S., Kemp, C.C.: Contactopt: Optimizing contact to improve grasps. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1471–1481 (2021)
23. Guzov, V., Sattler, T., Pons-Moll, G.: Visually plausible human-object interaction capture from wearable sensors. In: arXiv (May 2022)
24. Hamer, H., Gall, J., Weise, T., Van Gool, L.: An object-dependent hand pose prior from sparse training data. In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. pp. 671–678. IEEE (2010)
25. Hampali, S., Rad, M., Oberweger, M., Lepetit, V.: Honnotate: A method for 3d annotation of hand and object poses. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3196–3206 (2020)
26. Hasson, Y., Tekin, B., Bogo, F., Laptev, I., Pollefeys, M., Schmid, C.: Leveraging photometric consistency over time for sparsely supervised hand-object reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 571–580 (2020)
27. Hasson, Y., Varol, G., Laptev, I., Schmid, C.: Towards unconstrained joint hand-object reconstruction from rgb videos. arXiv preprint arXiv:2108.07044 (2021)
28. Hasson, Y., Varol, G., Tzionas, D., Kalevatykh, I., Black, M.J., Laptev, I., Schmid, C.: Learning joint reconstruction of hands and manipulated objects. In: CVPR (2019)
29. Hasson, Y., Varol, G., Tzionas, D., Kalevatykh, I., Black, M.J., Laptev, I., Schmid, C.: Learning joint reconstruction of hands and manipulated objects. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11807–11816 (2019)
30. Henter, G.E., Alexanderson, S., Beskow, J.: Moglow: Probabilistic and controllable motion synthesis using normalising flows. *ACM Transactions on Graphics (TOG)* **39**(6), 1–14 (2020)
31. HUANG, L., ZHANG, B., GUO, Z., XIAO, Y., CAO, Z., YUAN, J.: Survey on depth and rgb image-based 3d hand shape and pose estimation. *Virtual Reality & Intelligent Hardware* **3**(3), 207–234 (2021)
32. Jiang, H., Liu, S., Wang, J., Wang, X.: Hand-object contact consistency reasoning for human grasps generation. arXiv preprint arXiv:2104.03304 (2021)

33. Jiang, W., Kolotouros, N., Pavlakos, G., Zhou, X., Daniilidis, K.: Coherent reconstruction of multiple humans from a single image. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5579–5588 (2020)
34. Jiang, Z., Zhu, Y., Svetlik, M., Fang, K., Zhu, Y.: Synergies between affordance and geometry: 6-dof grasp detection via implicit representations. *Robotics: science and systems* (2021)
35. Karunratanakul, K., Yang, J., Zhang, Y., Black, M.J., Muandet, K., Tang, S.: Grasping field: Learning implicit representations for human grasps. In: 2020 International Conference on 3D Vision (3DV). pp. 333–344. IEEE (2020)
36. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: Bengio, Y., LeCun, Y. (eds.) 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14–16, 2014, Conference Track Proceedings (2014), <http://arxiv.org/abs/1312.6114>
37. Kocabas, M., Athanasiou, N., Black, M.J.: Vibe: Video inference for human body pose and shape estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5253–5263 (2020)
38. Kundu, J.N., Gor, M., Babu, R.V.: Bihmp-gan: Bidirectional 3d human motion prediction gan. In: Proceedings of the AAAI conference on artificial intelligence. vol. 33, pp. 8553–8560 (2019)
39. Kwon, T., Tekin, B., Stuhmer, J., Bogo, F., Pollefeys, M.: H2o: Two hands manipulating objects for first person interaction recognition. arXiv preprint arXiv:2104.11181 (2021)
40. León, B., Ulbrich, S., Diankov, R., Puche, G., Przybylski, M., Morales, A., Asfour, T., Moio, S., Bohg, J., Kuffner, J., et al.: Opengrasp: a toolkit for robot grasping simulation. In: International Conference on Simulation, Modeling, and Programming for Autonomous Robots. pp. 109–120. Springer (2010)
41. Li, J., Villegas, R., Ceylan, D., Yang, J., Kuang, Z., Li, H., Zhao, Y.: Task-generic hierarchical human motion prior using vaes. arXiv preprint arXiv:2106.04004 (2021)
42. Liu, C.K.: Dextrous manipulation from a grasping pose. In: ACM SIGGRAPH 2009 papers, pp. 1–6 (2009)
43. Luo, Z., Golestaneh, S.A., Kitani, K.M.: 3d human motion estimation via motion compression and refinement. In: Proceedings of the Asian Conference on Computer Vision (2020)
44. Malik, J., Abdelaziz, I., Elhayek, A., Shimada, S., Ali, S.A., Golyanik, V., Theobalt, C., Stricker, D.: Handvoxnet: Deep voxel-based network for 3d hand shape and pose estimation from a single depth map. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7113–7122 (2020)
45. Miller, A.T., Allen, P.K.: Graspit! a versatile simulator for robotic grasping. *IEEE Robotics & Automation Magazine* **11**(4), 110–122 (2004)
46. Mordatch, I., Popović, Z., Todorov, E.: Contact-invariant optimization for hand manipulation. In: Proceedings of the ACM SIGGRAPH/Eurographics symposium on computer animation. pp. 137–144 (2012)
47. Mueller, F., Davis, M., Bernard, F., Sotnychenko, O., Verschoor, M., Otaduy, M.A., Casas, D., Theobalt, C.: Real-time pose and shape reconstruction of two interacting hands with a single depth camera. *ACM Transactions on Graphics (TOG)* **38**(4), 1–13 (2019)
48. Ng, E., Ginosar, S., Darrell, T., Joo, H.: Body2hands: Learning to infer 3d hands from conversational gesture body dynamics. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11865–11874 (2021)

49. Oikonomidis, I., Kyriazis, N., Argyros, A.A.: Full dof tracking of a hand interacting with an object by modeling occlusions and physical constraints. In: 2011 International Conference on Computer Vision. pp. 2088–2095. IEEE (2011)
50. Ormoneit, D., Sidenbladh, H., Black, M.J., Hastie, T.: Learning and tracking cyclic human motion. *Advances in Neural Information Processing Systems* pp. 894–900 (2001)
51. Panteleris, P., Argyros, A.: Back to rgb: 3d tracking of hands and hand-object interactions based on short-baseline stereo. In: *Proceedings of the IEEE International Conference on Computer Vision Workshops*. pp. 575–584 (2017)
52. Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A., Tzionas, D., Black, M.J.: Expressive body capture: 3d hands, face, and body from a single image. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 10975–10985 (2019)
53. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 652–660 (2017)
54. Romero, J., Tzionas, D., Black, M.J.: Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)* **36**(6) (Nov 2017)
55. Sahbani, A., El-Khoury, S., Bidaud, P.: An overview of 3d object grasp synthesis algorithms. *Robotics and Autonomous Systems* **60**(3), 326–336 (2012)
56. Smith, B., Wu, C., Wen, H., Peluse, P., Sheikh, Y., Hodgins, J.K., Shiratori, T.: Constraining dense hand surface tracking with elasticity. *ACM Transactions on Graphics (TOG)* **39**(6), 1–14 (2020)
57. Sridhar, S., Mueller, F., Zollhöfer, M., Casas, D., Oulasvirta, A., Theobalt, C.: Real-time joint tracking of a hand manipulating an object from rgb-d input. In: *European Conference on Computer Vision*. pp. 294–310. Springer (2016)
58. Sridhar, S., Rhodin, H., Seidel, H.P., Oulasvirta, A., Theobalt, C.: Real-time hand tracking using a sum of anisotropic gaussians model. In: 2014 2nd International Conference on 3D Vision. vol. 1, pp. 319–326. IEEE (2014)
59. Starke, S., Zhang, H., Komura, T., Saito, J.: Neural state machine for character-scene interactions. *ACM Trans. Graph.* **38**(6), 209–1 (2019)
60. Taheri, O., Ghorbani, N., Black, M.J., Tzionas, D.: Grab: A dataset of whole-body human grasping of objects. In: *European Conference on Computer Vision*. pp. 581–600. Springer (2020)
61. Taylor, J., Bordeaux, L., Cashman, T., Corish, B., Keskin, C., Sharp, T., Soto, E., Sweeney, D., Valentin, J., Luff, B., et al.: Efficient and precise interactive hand tracking through joint, continuous optimization of pose and correspondences. *ACM Transactions on Graphics (TOG)* **35**(4), 1–12 (2016)
62. Taylor, J., Shotton, J., Sharp, T., Fitzgibbon, A.: The vitruvian manifold: Inferring dense correspondences for one-shot human pose estimation. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. pp. 103–110. IEEE (2012)
63. Taylor, J., Tankovich, V., Tang, D., Keskin, C., Kim, D., Davidson, P., Kowdle, A., Izadi, S.: Articulated distance fields for ultra-fast tracking of hands interacting. *ACM Transactions on Graphics (TOG)* **36**(6), 1–12 (2017)
64. Tiwari, G., Antic, D., Lenssen, J.E., Sarafianos, N., Tung, T., Pons-Moll, G.: Pose-ndf: Modeling human pose manifolds with neural distance fields. In: *European Conference on Computer Vision (ECCV)*. Springer (October 2022)
65. Urtasun, R., Fleet, D.J., Fua, P.: 3d people tracking with gaussian process dynamical models. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06). vol. 1, pp. 238–245. IEEE (2006)

66. Wang, Y., Min, J., Zhang, J., Liu, Y., Xu, F., Dai, Q., Chai, J.: Video-based hand manipulation capture through composite motion control. *ACM Transactions on Graphics (TOG)* **32**(4), 1–14 (2013)
67. Xie, X., Bhatnagar, B.L., Pons-Moll, G.: Chore: Contact, human and object reconstruction from a single rgb image. In: *European Conference on Computer Vision (ECCV)*. Springer (October 2022)
68. Yang, L., Zhan, X., Li, K., Xu, W., Li, J., Lu, C.: Cpf: Learning a contact potential field to model the hand-object interaction. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 11097–11106 (2021)
69. Ye, Y., Liu, C.K.: Synthesis of detailed hand manipulations using contact sampling. *ACM Transactions on Graphics (TOG)* **31**(4), 1–10 (2012)
70. Yi, H., Huang, C.H.P., Tzionas, D., Kocabas, M., Hassan, M., Tang, S., Thies, J., Black, M.J.: Human-aware object placement for visual environment reconstruction. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 3959–3970 (Jun 2022)
71. Zeng, A., Yang, L., Ju, X., Li, J., Wang, J., Xu, Q.: Smoothnet: A plug-and-play network for refining human poses in videos. In: *European Conference on Computer Vision*. Springer (2022)
72. Zhang, B., Wang, Y., Deng, X., Zhang, Y., Tan, P., Ma, C., Wang, H.: Interacting two-hand 3d pose and shape reconstruction from single color image. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 11354–11363 (2021)
73. Zhang, H., Bo, Z.H., Yong, J.H., Xu, F.: Interactionfusion: real-time reconstruction of hand poses and deformable objects in hand-object interactions. *ACM Transactions on Graphics (TOG)* **38**(4), 1–11 (2019)
74. Zhang, H., Zhou, Y., Tian, Y., Yong, J.H., Xu, F.: Single depth view based real-time reconstruction of hand-object interactions. *ACM Transactions on Graphics (TOG)* **40**(3), 1–12 (2021)
75. Zhang, H., Ye, Y., Shiratori, T., Komura, T.: Manipnet: neural manipulation synthesis with a hand-object spatial representation. *ACM Transactions on Graphics (TOG)* **40**(4), 1–14 (2021)
76. Zhang, S., Zhang, Y., Bogo, F., Pollefeys, M., Tang, S.: Learning motion priors for 4d human body capture in 3d scenes. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 11343–11353 (2021)
77. Zhang, X., Bhatnagar, B.L., Guzov, V., Starke, S., Pons-Moll, G.: Couch: Towards controllable human-chair interactions. In: *European Conference on Computer Vision (ECCV)*. Springer (October 2022)
78. Zhao, R., Su, H., Ji, Q.: Bayesian adversarial human motion synthesis. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 6225–6234 (2020)
79. Zhao, W., Zhang, J., Min, J., Chai, J.: Robust realtime physics-based motion control for human grasping. *ACM Transactions on Graphics (TOG)* **32**(6), 1–12 (2013)
80. Zhao, Z., Wang, T., Xia, S., Wang, Y.: Hand-3d-studio: A new multi-view system for 3d hand reconstruction. In: *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. pp. 2478–2482. IEEE (2020)
81. Zhu, T., Wu, R., Lin, X., Sun, Y.: Toward human-like grasp: Dexterous grasping via semantic representation of object-hand. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 15741–15751 (2021)
82. Zimmermann, C., Ceylan, D., Yang, J., Russell, B., Argus, M., Brox, T.: Freihand: A dataset for markerless capture of hand pose and shape from single rgb images.

In: Proceedings of the IEEE/CVF International Conference on Computer Vision.
pp. 813–822 (2019)