# ArtiBoost: Boosting Articulated 3D Hand-Object Pose Estimation via Online Exploration and Synthesis

[1]Kailin Li[*], [1,2]Lixin Yang[*], [1]Xinyu Zhan, [1]Jun Lv, [1,2]Wenqiang Xu, [1]Jiefeng Li, [1,2]Cewu Lu[†]
[1]Shanghai Jiao Tong University, China   [2]Shanghai Qi Zhi Institute, China
{kailinli, siriusyang, kelvin34501, LyuJune_SJTU, vinjohn, ljf_likit, lucewu}@sjtu.edu.cn

## Abstract

*Estimating the articulated 3D hand-object pose from a single RGB image is a highly ambiguous and challenging problem, requiring large-scale datasets that contain diverse hand poses, object types, and camera viewpoints. Most real-world datasets lack these diversities. In contrast, data synthesis can easily ensure those diversities separately. However, constructing both valid and diverse hand-object interactions and efficiently learning from the vast synthetic data is still challenging. To address the above issues, we propose ArtiBoost, a lightweight online data enhancement method. ArtiBoost can cover diverse hand-object poses and camera viewpoints through sampling in a Composited hand-object Configuration and Viewpoint space (CCV-space) and can adaptively enrich the current hard-discernable items by loss-feedback and sample re-weighting. ArtiBoost alternatively performs data exploration and synthesis within a learning pipeline, and those synthetic data are blended into real-world source data for training. We apply ArtiBoost on a simple learning baseline network and witness the performance boost on several hand-object benchmarks. Our models and code are available at https://github.com/lixiny/ArtiBoost.*

## 1. Introduction

Articulated bodies, such as the human hand, body, and linkage mechanism, can be observed every day in our life. Their joints, links, and movable parts depict the functionality of the articulation body. Extracting their transient configuration from image or video sequence, which is often referred to as *Pose Estimation* [8, 37, 38], can benefit many downstream tasks in robotics and augment reality. Pose estimation for multi-body articulations is especially challenging as it suffers from severe self- or mutual occlusion. In
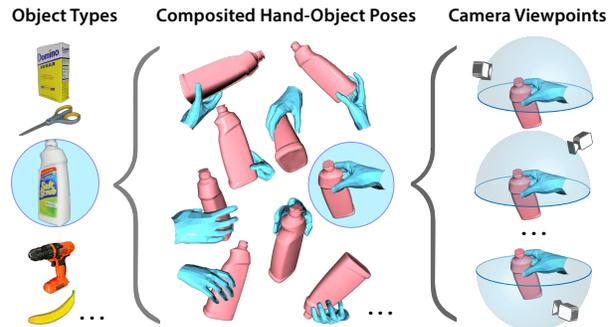


Figure 1. **An intuitive illustration of the CCV-space.**

this work, we paid attention to a certain type of multi-body articulations – composited hand and object poses during their interaction [5, 12, 19, 21, 24–26, 30, 32, 43, 69]. Hands are the primary means by which humans manipulate objects in the real-world, and the hand-object pose estimation (HOPE) task holds great potential for understanding human behavior [14, 18, 36, 46, 61].

As the degrees of freedom (DoF) grows, the proper amount of data to cover the pose distribution has grown exponentially. More than the most common articulation bodies, the human hand has 16 joints and approximately 21 DoF. Preparing such diverse training data for the HOPE task can be very challenging. The real-world recording and annotation methods [3, 6, 14, 21] tend to hinder the pose diversity. For example, the multi-view-based approaches [3, 21] require the subject to maintain a static grasping pose in a video sequence. As a result, their recording process is inefficient, and their pose diversity is insufficient. In contrast, data synthesis is efficient and annotation-free, and has been widely adopted in single-body N-D pose estimation [9, 15, 42, 63, 65, 72]. However, these methods are ineligible for multi-body articulation, in which the poses are restricted by mutual contact and obstruction. Data synthesis for HOPE tasks requires us to simulate virtual hand-object interaction (grasp) that mimics the underlying pose distribution of their real-world counterparts. Conventional method either manually articulated [10, 51, 52] the hand model for grasp, or relegated [2, 26, 34] grasp synthesis to an off-the-

---

shelf grasping simulator: GraspIt [48]. However, the manual methods are difficult to scale their data in large amounts, and the simulation method also sacrificed the diversity of hand poses. GraspIt optimizes for hand-crafted grasp metrics [13] that do not reflect the pose distribution of a 21-DoF dexterous hand. Besides, even with a vast amount of synthetic grasps, not every hand-object configuration is helpful for training. For example, similar configurations may have already been observed multiple times, and those easily discernable samples may have a frequent appearance. Hence, offline data synthesis, without repeatedly communicating with the model during training, is still considered inefficient for a learning task.

To address the above issues, we propose an online data enhancement method **ArtiBoost**, to effectively **boost** the **arti**culated hand-object pose estimation via two alternative steps, namely exploration and synthesis. First, to describe the observation of hand-object interaction, we design a three-dimensional discrete space: **C**omposited hand-object **C**onfiguration and **V**iewpoint space (CCV-space) where object types, hand pose, and viewpoint are its components. Second, to construct valid and diverse hand-object poses in CCV-space, we design a fitting-based grasp synthesis method that exploits the contact constraints [69] between hand and object vertices to simulate MANO hand [53] grasping a given object. After the CCV-space is established, we next describe how ArtiBoost enhances the HOPE tasks.

At the exploration step, ArtiBoost explores the CCV-space and samples different hand-object-viewpoint triplets from it. Then at the synthesis step, the hand and object in the triplet will be rendered on the image from the viewpoint in the triplet. These synthetic images are mixed with the real-world source images in batches to train the HOPE model. Later, the training losses are fed back to the exploration step and guide it to re-weight the current hard-discernable triplets for the next round of sampling. With such communication in the training loop, ArtiBoost can adaptively adjust its sampling weights to select more hard-discernable data for the current HOPE model. As the HOPE model becomes powerful, it can also continuously promote the current ArtiBoost to evolve. ArtiBoost is model-agnostic, which means it can be plugged into any modern CNN architecture. In this paper, we plug ArtiBoost into a simple classification-based (*e.g.* [57]) and regression-based (*e.g.* [1]) pose estimation model to show its efficacy. For evaluation, we report those models' performance on two challenging HOPE benchmarks: HO3D (v1-v3) [20–22] and DexYCB [6]. Without whistles and bells, those simple baseline models can outperform the results of previous state-of-the-arts.

In this paper, we propose to boost the performance of HOPE task by enhancing the diversity of underlying poses distribution in the training data. We summarize our contributions as follows. (1) To describe the composited hand-object-viewpoint poses distribution, we design the CCV-space. (2) To overcome the scarcity of such poses in the previous dataset, we design a contact-guided grasp synthesis method and simulate both valid and diverse hand-object poses to fill the CCV-space. (3) To help HOPE model efficiently fit the underlying poses distribution, we parallelize the data synthesis with the learning pipeline, leverage the training feedback, and adopt a sample re-weighting strategy. Finally, We conduct extensive experiments to validate our technical contributions (Sec. 4.2, 4.3) and reveal the potential applications (Sec. 4.4).

## 2. Related Work

**Hand-Object Pose Estimation.** As some HOPE tasks are closely related to hand pose estimation (HPE) tasks [64, 66], we firstly review several HPE methods. According to its output form, single RGB-based 3D HPE can be categorized into three types: image-to-pose (I2P) [68], image-to-geometry (I2G) [8, 60], and hybrid [71]. While the I2P only focuses on the joints' pose only, the I2G focus on recovering the full hand geometry (pose and shape). Meanwhile, recent works [37, 67, 71] showed that I2G could be hybridized to I2P through neural inverse kinematics (IK). Second, we explore several HOPE methods. Regarding the learning-based methods, some aimed to predict the hand-object poses in a unified model [24, 43], while the others focused on recovering hand-object interaction based on contact modeling [19, 69]. As for the learning-free methods, Hasson *et al*. [25] and Cao *et al*. [5] proposed to aggregate the visual cues from object detection, HPE, and instance segmentation to acquire the optimal hand-object configuration. This paper adopts two simple learning baseline networks of two paradigms: classification (joints as 3D heatmap) and regression (joints from pose and shape predictions). We show that with ArtiBoost, even simple baseline networks can outperform previous sophisticated CNN designs.

**Data Synthesis for Pose Estimation.** Using synthetic data to increase pose variants has been widely adopted in single-body N-D pose estimation tasks, such as human pose [7, 63, 65], hand pose [9, 15, 72], 6D object pose [33, 58, 62], and 7D articulated object pose [38, 42], *etc*. These methods utilized the kinematic model of the articulated bodies, drive constrained joints' motion, and rendered the current model from different viewpoints. Unlike single-body pose estimation, Data synthesis for HOPE (multi-body) task not only needs to consider the joints limit, but also to obey the physical constraints brought from mutual contact and obstruction. In terms of composited hand-object pose synthesis, there are also three genres included in the literature. The manual labeling methods [10, 51, 52] articulated a hand model to achieve grasp; The metric-based methods [2, 26, 34] leveraged grasp simulator [48] to simulate
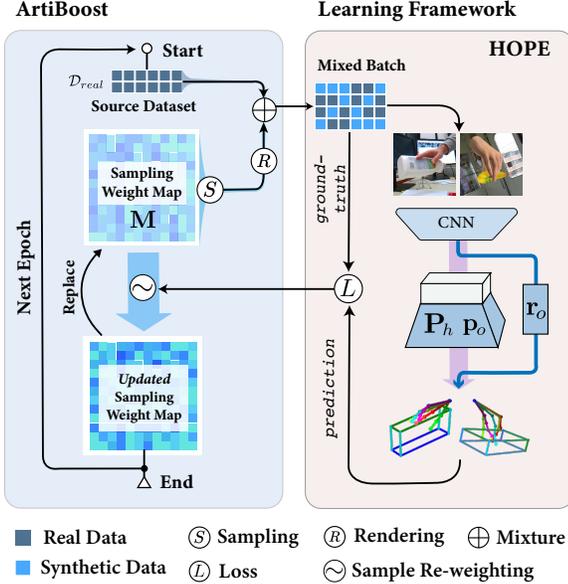
Figure 2. **Illustration of the integrated pipeline.** ArtiBoost can be plugged into an arbitrary HOPE framework by modifying the current data loader.

hand poses subjected to physical grasping metrics; The data-driven methods [31, 32, 59] trained conditional VAE that generates new grasps. This paper presents an automatically contact-guided optimization method to construct valid and diverse poses of hand-object interaction.

**Exploration for Hard Examples,** also called "hard example/negative mining", has been proven effective for various computer vision tasks, such as object detection [40, 55], person re-id [28], head pose estimation [35], face recognition [54], and deep metric learning [16, 56]. Generally speaking, the basic ideology of hard examples mining is that if a prediction of a certain data sample exhibits a large error under a certain metric, then this data sample is not properly learned by the learning algorithm. By adding such data samples to the training batch can help the learning converge faster. Recent work exploited a generative adversarial model [17] to acquire harder samples based on error feedback training. However, it must pay non-trivial efforts on the adversarial part to ensure samples' validity. In this paper, we adopt a simple yet effective sample re-weighting strategy that adaptively selects those hard-discernable training triplets (hand-object-viewpoint) for training.

## 3. Method

**Overview.** This section describes the exploration and synthesis step in ArtiBoost and elaborates the learning framework for the HOPE task. Inside the explorations step, we present the composited configurations and viewpoints space (CCV-space), the key component of ArtiBoost.

**Problem Definition.** Given an input image $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$

that observes a single hand interacting with a certain object, HOPE aims to learn a certain neural network that predict the 3D hand joint locations: $\mathbf{P}_h = \{\mathbf{p}_j\}_{j=1}^J$, object centroid locations: $\mathbf{p}_o$ and object rotation: $\mathbf{r}_o \in \mathfrak{so}(3)$, where $\mathbf{p}_j, \mathbf{p}_o \in \mathbb{R}^3$, $J = 21$ and the $H \times W$ is the resolution.

To train the neural network, we shall firstly prepare a real-world source dataset: $\mathcal{D}_{real}$. ArtiBoost is employed along with the $\mathcal{D}_{real}$. During the training process, ArtiBoost iteratively samples (without replacement) hand-object-viewpoint triplets from the CCV-space: $\mathcal{C}$ based on a weight map: $\mathbf{M}$. Each entry of $\mathbf{M}$ corresponds to a certain triplet in $\mathcal{C}$, and the value in that entry corresponds to the sampling weight of the triplet. Meanwhile, those selected triplets are rendered as a batch of synthetic images at the synthesis step. The synthetic images are mixed with source images from $\mathcal{D}_{real}$. After that, the mixed batch is fed to the HOPE learning framework to complete a forward and backward propagation. When an epoch of training has finished, ArtiBoost performs the sample re-weighting in $\mathbf{M}$ based on the loss value and waits for the next round of training. The whole pipeline is illustrated in Fig. 2.

### 3.1. Online Exploration in CCV-Space

**The Composited Configuration & Viewpoint Space.** HOPE problem commonly involves a certain interacting hand-object configuration that is observed by a certain viewpoint. The input domain of HOPE can thus be narrowed down to three main dimensions: object type, hand pose, and viewpoint direction. To note, the dimension of object type and hand pose are not independent of each other. Given a certain object model, the hand pose that interacts with it depends on the geometry of the model. As shown in Fig. 1, we define the discrete representation of the input domain as the CCV-space: $\mathcal{S} = \{(n_o, n_p, n_v) \in \mathbb{N}_+^3 \mid n_o \leq N_o, \ n_p \leq N_p, n_v \leq N_v\}$, where the $N_o$, $N_p$ and $N_v$ is the number of object types, discrete poses and viewpoints, respectively. The $(i, j, k)$ item in $\mathcal{S}$ stands for the scenario that the interaction between the $i$-th object and the $j$-th hand pose is observed at the $k$-th camera viewpoint. Next, we will sequentially present the components in CCV-space, namely: hand configuration space (**C-space**), composited hand-object configuration space (**CC-space**) and viewpoint space (**V-space**).

**C-Space of Valid Hand Pose.** To represent hand, we employ a parametric skinning hand model, MANO [53] which drives a deformable hand mesh with 16 joints rotations $\boldsymbol{\theta} \in \mathbb{R}^{16 \times 3}$ and shape parameters $\boldsymbol{\beta} \in \mathbb{R}^{10}$. Given the axis-angle forms of rotation, MANO has 48 DoFs that exceed the DoFs allowed by a valid hand pose [39]. Fitting or interpolation on the 48 DoFs rotations may encounter abnormal hand pose that is unhealthy for training the HOPE network. Besides, the original MANO's coordinate system is not coaxial with the direction of the hand's kinematic tree,
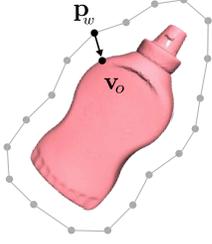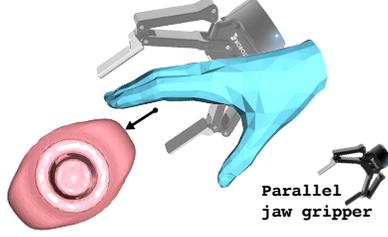
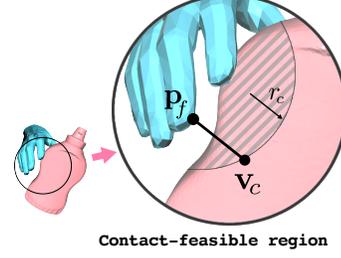Figure 3. Offset surface.     Figure 4. Pre-grasping hand pose.     Figure 5. Fingertip and its contact point.

so that the rotation axis is coupled with at least two of the MANO's orthogonal axes. This property makes the pose interpolation more difficult. In the paper, we employ an axis adaptation based on the *twist-splay-bend* coordinate system that was initially proposed by [69]. It enables us to describe the hand pose at each joint as the rotation angle along with one of the specified coordinate axis (*e.g. bend*) at this joint. With the axis adaptation, we design three protocols for describing the C-space of valid hand pose:

**i).** According to hand anatomy, all the non-metacarpal joints along the hand's kinematic tree can only have the bending pose. And the five metacarpal joints can only have a combined bending and splaying poses. Any twisting along the pointing direction or splaying at non-metacarpal joints is prohibited.

**ii).** For each of the five fingers, the bending poses on their proximal and distal joint are linked and dependent. The bending poses at the five metacarpal joints are independent of other joints.

**iii).** The pose of each finger is independent of each other as long as it does not conflict with the protocol **i)** & **ii)**.

Based on these protocols, the total DoFs in the hand's C-space is 21 (one splaying and two independent bending DoFs for each of the five fingers, plus 6 DoFs at the wrist). Hence, we can describe the whole hand pose by describing the angles of 15 joints' rotation along its specified axis: $\theta_i^{bend}$ or $\theta_i^{splay}$, as well as the wrist pose: $\boldsymbol{\xi}_w \in \mathfrak{se}(3)$. These protocols not only guarantee the diverse and valid hand poses but also ensure these poses are visually plausible. After constructing a valid C-space, we now move on to the composited space that describes hand-object interaction.

**Composited C-Space of Hand-Object Interaction.** We firstly define the hand-object "interaction" as the scenario that satisfies the following requirements: **i)** The thumb and at least one or more fingers should be in contact with the object surface [59]; **ii)** The hand and object model should not intersecting with each other; As the hand's interacting pose is highly dependent on its approaching direction and the object's geometry, only a small portion of poses in hand C-space is valid for interaction purposes. Hence instead of searching for valid interaction poses in the entire hand C-space, we turn to explore a discrete space of predefined

hand-object poses. We call it composited C-space (CC-space). The goal of CC-space is to increase the diversity of the interacting poses for a given object. For this purpose, we leverage the contact constraints to optimize grasps. This method is divided into three steps.

● **1).** First, given an object model, we construct an offset surface outside of its original surface and uniformly sample $N_w$ points on the offset surface (Fig. 3). These points will control the wrist position. For each point: $\mathbf{p}_w$, we query its closet vertex on object surface: $\mathbf{v}_o$. The vector: $(\mathbf{v}_o - \mathbf{p}_w)$ represents the approaching direction from hand to object. Then, we construct a prehensile hand pose that mimics the rest state of a parallel jaw gripper (Fig. 4). This pre-grasping hand is placed at each $\mathbf{p}_w$ as interaction initiation. The $\mathbf{p}_w$'s approaching direction controls the movement of the wrist. Based on a certain $\mathbf{p}_w$, we define the region between $\mathbf{v}_o$ and the farthest vertex that the finger can reach as the object's contact-feasible region.

● **2).** Second, to generate the interacting pose, we fit the fingertips to contact point $\mathbf{v}_c$ chosen from contact-feasible region (Fig. 5). For each $\mathbf{p}_w$, we generate $N_p^w$ interacting hand poses. To increase diversity, we employ several randomness during fitting: **i)** We randomly select thumb and $N$ $(1 \le N \le 4)$ other fingers on hand. Only the selected fingers will participate in fitting. **ii)** For each selected fingers, we set a random minimal reaching radius $r_c$ inside the contact-feasible region. Now the selected fingertip must reach for a contact point $\mathbf{v}_c$ that satisfies: **i)** $\|\mathbf{v}_c - \mathbf{p}_w\|_2 \ge r_c$; **ii)** $\min_{\mathbf{v}_c}\|\mathbf{v}_c - \mathbf{p}_f\|_2$, where the $\mathbf{p}_f$ is the selected fingertip point. After each selected fingertip is paired with a certain contact point, we will initiate the fitting process.

● **3).** During fitting, we adopt the anchor-based hand model and contact-based cost function defined in CPF [69]. The fitting process aims to minimize the cost brought from unattached fingertips and contact points: the unattached anchor on a fingertip: $\mathbf{p}_f$ will be attracted to its corresponded contact point $\mathbf{v}_c$ on the object, while the intersected anchors will be pushed out. The fitting of hand poses is only performed on the predefined DoFs in the hand's C-space. Thus valid hand pose can be guaranteed.

**Viewpoint Space.** For the camera viewpoint, we uniformly sample $N_v$ viewpoints direction $\mathbf{n}_v$ according to the sphere

4

sampling strategy [45]:

$$\mathbf{n}_v = (\sqrt{1-u^2}\cos(\phi), \sqrt{1-u^2}\sin(\phi), u)^\mathsf{T} \quad (1)$$

where $u \sim \mathcal{U}[-1,1]$ and $\phi \sim \mathcal{U}[0,2\pi]$, where $\mathcal{U}$ stands for uniform distribution.

**CCV-space Implementation.** We iteratively fit 300 different interacting poses for each object. Since optimization may result in local minima, we manually discard the poses that **i)** exhibit severe inter-penetration between hand and object; **ii)** form an unnatural grasping or interaction. Then, we select up to 100 interacting hand poses per given object. Though the optimization can potentially generate unlimited interacting poses, we find 100 poses (per object) are sufficient to boost current HOPE tasks. In the viewpoints space, we choose $N_u = 12$ and $N_\phi = 24$, which comprise $N_v = N_u \times N_\phi = 288$ different viewpoints. The total number of different (hand-object-viewpoint) triplets depends on the benchmarking dataset. For example, the triplets catered for DexYCB dataset (containing 20 YCB objects) is $N_o \times N_p \times N_v = 20 \times 100 \times 288 = 576,000$.

**Weight-guided Sampling Strategy.** In literature, the uniform sampling strategy was widely adopted by synthetic dataset [26, 49, 72]. However, not every sample in the CCV-space contributes equally to the network. Since we hope that those hardly discernable samples shall have a higher frequency of occurrence, we construct a weight map $\mathbf{M} \in \mathbb{R}^{N_o \times N_p \times N_v}$ to guide the sampling in the exploration step. In $\mathbf{M}$, each element $w_i$ stands for the sampling weight of the corresponding item in CCV-space. The probability $p_i$ of a certain item that would be sampled is $p_i = w_i / \sum w_j$. We then draws $N_{syn}$ samples from the multinomial distribution $\{p_i \mid p_i = w_i / \sum w_j; \ w_i, w_j \in \mathbf{M}\}$. Based on this strategy, we shall increase the weights for those hardly discernable samples in $\mathbf{M}$ while decreasing the weights for those who are already easy to discern, when we get the feedback from the loss value.

**Sample Re-weighting.** After an epoch of training has finished, ArtiBoost chooses those hard discernable items and re-weights their sampling weights. We inspect a percentile-based re-weighting strategy. During the re-weighting phase, each synthetic sample will be assigned a weight update. These updates are multiplied by the original sampling weight in $\mathbf{M}$. Intuitively, we want those hard discernable samples to have high weight. In the percentile-based re-weighting strategy, we calculate the weight update based on the percentile of the samples' Mean Per Joint Position Error (MPJPE) among the whole epoch of synthetic samples. For the $i$-th sample, given by the MPJPE $e_i$ and its percentile $q_i = \frac{e_{\max} - e_i}{e_{\max} - e_{\min}}$, the weight update are calculated from a simple reciprocal heuristic: $\delta w_i = \frac{1}{q_i + 0.5}$. If the sample $i$ has the maximum MPJPE $e_{\max}$ among the synthetic samples in current epoch, its original sampling weight in $\mathbf{M}$ will be multiplied by a maximum update factor $\delta w_i = 2$. If the $i$ has the minimum MPJPE $e_{\min}$, its update factor would be $\delta w_i = 2/3$. We also clamp the updated $\mathbf{M}$ by a upper bound 2.0 and lower bound 0.1 to avoid over imbalance.

## 3.2. Online Synthesis for HOPE task

During the training process, we synthesize the sampled hand-object-viewpoint triplets to RGB images. This synthesis process is task-oriented, as the adaptive sampling decides its composition to cater to the downstream task. Here, we describe the features in the online synthesis step.

**Disturbance on the Triplets.** To increase the variance in the pose distribution and thus to improve the network's generalization ability, we add disturbance on the hand poses and viewpoint directions before rendering images.

• *For the hand poses*, we relieve the restriction in protocol **ii)** in hand C-space, in which the bending angles of distal and proximal joints on each finger are now independent in terms of disturbance. Then, we add a Gaussian disturbance $\mathcal{N}(0, \sigma_1^2)$ on each of the 15 bending angles. Second, for the disturbance on splaying angles, we add a $\mathcal{N}(0, \sigma_2^2)$ on the five metacarpal joints. We empirically set $\sigma_1 = 3$ and $\sigma_2 = 1.5$ degree. To note, this disturbance still subject to the restrictions in protocol **i)** and **iii)**, which ensure a valid and prehensile hand configuration. However, these disturbances may cause the inter-penetration between the hand and object models. Hence, we further process the disturbed hand-object pose through the *RefineNet* module in GrabNet [59], which has the effect of mitigating inter-penetration. Apart from hand pose, the shape of hand also impacts the interaction with the object. Hence we sample the random MANO shape parameters ($\beta \in \mathbb{R}^{10}$) from the distribution of $\mathcal{N}(0, 0.5)$ to formulate the final hand model.

• *For the viewpoint directions*, we add three disturbances: $\mathcal{U}(-\delta u, +\delta u)$, $\mathcal{U}(-\delta\phi, +\delta\phi)$, and $\mathcal{U}(0, 2\pi)$ on the elevation distance $u$, azimuth angle $\phi$ and the camera in-plane rotation, respectively. In all experiments, we empirically set $\delta u = 0.05$ and $\delta\phi = 7.5$ degree.

**Skin Tone & Textures.** We adopt a state-of-the-art hand skin tone & texture model: HTML [50] for realistic appearance on the rendered images. HTML represents the hand's skin color & texture of as continuous parameters in a PCA space. Before the refined hand model enters the rendering pipeline, we randomly assign it an HTML texture map. We also found that discarding the shadow removal operation in HTML produces more visually plausible images.

**Rendering.** We employ the off-the-shelf rendering software: PyRender [47] inside ArtiBoost. The interactive hand and object are the foreground, and the images in COCO [41] dataset are the background. The pipeline can support rendering of the photorealistic hand and object textured meshes onto a $224 \times 224$ canvas at 120 FPS per graphic card (Titan X), sufficient for us to parallelize the rendering with

Figure 6. The rendered images during online synthesis.

training. The synthetic data is the combination of four customized items, namely hand-object-viewpoint triplet, background, skin tone, and texture. We show several synthetic images in Fig. 6.

### 3.3. Learning Framework

We adopt two simple baseline network: one is classification-based (*Clas*) and the other regression-based (*Reg*). We employ ResNet-34 [27] as the backbone in both of them. In *Clas*, we use two de-convolution layers to generates 22 3D-heatmaps that indicate the location of 22 joints (21 hand joints and one object centroid) as likelihood. The 22 3D-heatmaps are defined in a restricted *uvd* space, where *uv* is the pixel coordinates, and *d* is the a wrist-relative depth value. Then, we employ a soft-argmax operator to convert the 3D-heatmaps into joints' *uvd* coordinates. Finally, we transform the joints' *uvd* coordinates into its 3D locations: $\mathbf{P}_h, \mathbf{p}_o$ in camera space by camera intrinsic: $\mathbf{K}$ and wrist location: $\mathbf{p}_w$. In *Reg*, we use multi-layer perceptron (MLP) to regress the MANO parameters: $\boldsymbol{\theta}$, $\boldsymbol{\beta}$ and object centroid: ${}^w\mathbf{p}_o$ w.r.t. the wrist. Then we transfer the $\boldsymbol{\theta}$, $\boldsymbol{\beta}$ to the wrist-relative hand joints: ${}^w\mathbf{P}_h$ by the MANO model. Finally the ${}^w\mathbf{P}_h$ and ${}^w\mathbf{p}_o$ are translated into the camera space by adding a known wrist location: $\mathbf{p}_w$. Both *Clas* and *Reg* adopt another MLP branch to predict object rotation: $\mathbf{r}_o$. Detailed implementations are provided in **Appx**.

**Loss Function.** The loss function to train the HOPE network consists of four terms.

First, we penalize the error of total 22 joints location (21 hand joints and 1 object centroid) in form of $\ell 2$ distance:

$$\mathcal{L}_{loc} = \frac{1}{22} \sum_{i=1}^{22} \left\| \mathbf{p}_i - \hat{\mathbf{p}}_i \right\|_2^2 \qquad (2)$$

where $\hat{\mathbf{p}}_i$ denotes the ground-truth joint location.

Second, we penalize the error of object rotation in form of $\ell 2$ distance at the eight tightest bounding box corners:

$$\mathcal{L}_{cor} = \frac{1}{8} \sum_{i=1}^{8} \left\| \exp(\mathbf{r}_o) * \bar{\mathbf{c}}_i - \hat{\mathbf{c}}_i \right\|_2^2 \qquad (3)$$

where the $\bar{\mathbf{c}}_i$ and $\hat{\mathbf{c}}_i$ are the object's corners in canonical view and camera view. $\exp(\mathbf{r}_o)$ is the predicted rotation.

Third, we adopt the ordinal relation loss to correct the 2D-3D misalignment. $\mathcal{L}_{ord}$ inspects the joint-level depth relation inside a pair of joints: one from the 21 hand joints and the other from the 8 object corners. We penalize the case if the predicted depth relation between the $i$-th hand joint: $\mathbf{p_i}$ and the $j$-th corner $\mathbf{c_j}$ is misaligned with its ground-truth

relation: $\mathbb{1}_{i,j}^{ord}$. The $\mathcal{L}_{ord}$ is formulated as:

$$\mathcal{L}_{ord} = \sum_{j=1}^{8} \sum_{i=1}^{J=21} \mathbb{1}_{i,j}^{ord} * \left| (\mathbf{p}_i - \mathbf{c}_j) \cdot \mathbf{n}_\perp \right| \qquad (4)$$

where the $\mathbf{n}_\perp$ is the viewpoint direction.

Fourth, we borrow a symmetry-aware object corner loss: $\mathcal{L}_{sym}$ from Hampali *et al.* [23]:

$$\mathcal{L}_{sym} = \min_{\mathbf{R} \in \mathcal{S}} \frac{1}{8} \sum_{i=1}^{8} \left\| \exp(\mathbf{r}_o) * \bar{\mathbf{c}}_i - \exp(\hat{\mathbf{r}}_o)\mathbf{R} * \bar{\mathbf{c}}_i \right\|_2^2 \quad (5)$$

where $\exp(\hat{\mathbf{r}}_o)$ denotes object's ground-truth rotation matrix. Given an object, the set $\mathcal{S}$ contains all the valid rotation matrices based on the object's predefined symmetry axes.

The overall loss is a weighted sum of the four terms:

$$\mathcal{L}_{HOPE} = \mathcal{L}_{loc} + \lambda_1 \mathcal{L}_{cor} + \lambda_2 \mathcal{L}_{ord} + \lambda_3 \mathcal{L}_{sym} \qquad (6)$$

where the $\lambda_{1\sim3}$ are the hyper-parameters.

## 4. Experiment and Result

### 4.1. Dataset and Metrics

**Dataset.** We evaluate our methods on three hand-object dataset: **FHAB** [14], **HO3D** [21] and **DexYCB** [6]. FHAB contains 20K samples of hand in manipulation with objects. We follow the "action" split as in Tekin *et al.* [61], which contains 10,503 training and 10,998 testing samples. The FHAB dataset only contains a few numbers of hand poses and viewpoints. We find its training set is adequate for the neural network. Thus we only use FHAB to verify the feasibility of the learning framework. HO3D is a dataset that contains a large number of images of hand-object interactions. Evaluation of the HO3D testing set is conducted at an online server. We also report our results on the latest **HO3Dv3** [22], which is released with different training/testing split. DexYCB contains 582K image frames of grasping on 20 YCB objects. We only evaluate the right-hand pose using the official "S0" split and filter out the frames that the minimum hand-object distance is large than 5 $cm$ to make sure a plausible hand-object interaction would appear.

**Metrics.** For the hand pose, we report the mean per joint position error (**MPJPE**) in the wrist-aligned coordinates system. For the object pose, there are two standard metrics in literature: mean per corners position error (**MPCPE**) and maximum symmetry-aware surface distance (**MSSD**). The former MPCPE directly measures the unique pose of the object. However, since some objects are symmetrical or revolutionary invariant, and since the objects are often severely occluded by hand, direct measuring objects' absolute and unique pose is sometimes less reasonable. MSSD measures the difference between the current object pose to its closest counterpart in all its rotation invariants. In this paper, we report the object's MPCPE and MSSD within dif-

| Method | MPJPE | MPCPE |
|---|---|---|
| Hasson *et al.* [24] | 11.33 | 28.42 |
| Our *Clas* w/o $\mathcal{L}_{ord}$ | 8.71 | **18.64** |
| Our *Clas* | **8.60** | 19.45 |

Table 1. Comparisons with SOTA on **FHAB** dataset (errors are reported in $mm$). The comparisons are made in the wrist-aligned coordinates system.

| Method | MPJPE | MPCPE |
|---|---|---|
| Hasson *et al.* [24] | 3.69 | 12.38 |
| Liu *et al.* [43] | 2.93 | - |
| Our *Reg* | 3.53 | 7.38 |
| Our *Reg* + **Arti** | **3.17** | **5.87** |
| Our *Clas* | 3.06 | 7.24 |
| Our *Clas* + **Arti** | **2.64** | **5.16** |

Table 2. Comparisons ($cm$) with SOTA on **HO3D** dataset.

| Method | MPJPE | MSSD *mustard bottle* | MSSD *bleach cleanser* | MSSD *potted meat can* |
|---|---|---|---|---|
| Hampali *et al.* [23] | 2.57 | 4.41 | 6.03 | 9.08 |
| Our *Clas* sym | 3.10 | 4.07 | 6.56 | 8.70 |
| Our *Clas* sym + **Arti** | **2.53** | **3.14** | **5.72** | **6.36** |

Table 3. Comparison ($cm$) with Transformer-based SOTA on **HO3D** using symmetry-aware loss $\mathcal{L}_{sym}$. We use the same symmetry axes as described in [23].

| Method | MPJPE | MSSD ★ *power drill* | MSSD ★ *cracker box* | MSSD ★ *scissors* | MSSD ★ *bleach cleanser* |
|---|---|---|---|---|---|
| Our *Clas* sym | 13.00 | 74.95 | 63.68 | 88.10 | 91.66 |
| Our *Clas* sym + **Arti** | **12.80** | **52.70** | **46.13** | **66.52** | **72.31** |

Table 4. Our results ($mm$) on **DexYCB**. ★ We only list the MSSD score of 4 objects. The full table can be found in **Appx**.

| Method | MPJPE | MPCPE | MSSD *mustard bottle* | MSSD *bleach cleanser* | MSSD *potted meat can* |
|---|---|---|---|---|---|
| Our *Clas* | 2.94 | 7.53 | 6.88 | 5.56 | 7.63 |
| Our *Clas* + **Arti** | 2.50 | 5.88 | 3.79 | **4.99** | 6.21 |
| Our *Clas* sym | 2.98 | - | 3.73 | 6.39 | 7.28 |
| Our *Clas* sym + **Arti** | **2.34** | - | **2.66** | 5.23 | **5.82** |

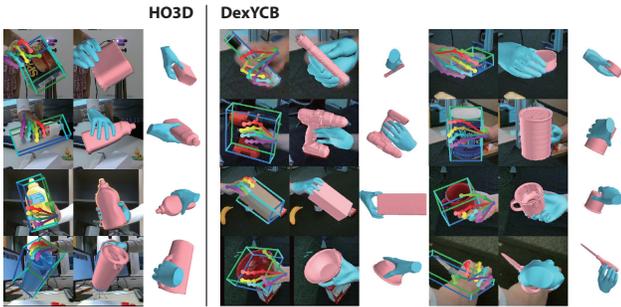Table 5. Our results ($cm$) on **HO3D v3**.



Figure 7. Qualitative results on the **HO3D** and **DexYCB** dataset.

ferent training schemes. When reporting MPCPE, we train the network with $\lambda_1 = \lambda_2 = 1$ and $\lambda_3 = 0$. When taking objects' symmetricity into account, we train the network with $\lambda_1 = \lambda_2 = 0$ and $\lambda_3 = 1$. We call the later network symmetry model (abbr. sym) and report its MSSD following the BOP challenge protocol [29]. The definition of YCB objects' symmetry axes can be found in **Appx**.

## 4.2. HOPE Network Performance

**Qualitative Results.** Our qualitative results on HO3D and DexYCB testing sets are shown in Fig. 7. We draw the predicted hand joints and object corners as their 2D projections (1st col). We also adopt a pretrained IK network [44] that maps the hand joints to hand mesh surfaces for visualization. We draw the full hand-object geometry in camera view (2nd col) and another viewpoint (3rd col). More qualitative results are provided in **Appx**.

**Comparison with State-of-the-Art.** In Tab. 1 we compare our *Clas* with the previous SOTA [24] on FHAB to justify the ordinal relation loss $\mathcal{L}_{ord}$. To note, [24] regressed hand and object poses in camera space. For fair comparison, we align their results in wrist-relative coordinates system.

In the subsequent tables, we use "**+ Arti**" to denote a certain network is trained by the original dataset's training split **plus** the synthetic data brought from **Arti**Boost's exploration and synthesis. In Tab. 2, we show that our ArtiBoost enhances both *Reg* and *Clas* performance on HO3D dataset. We obtain **28%** and **10%** MPJPE improvement compared with the previous SOTA [24] and [43], respectively. For fair compassion, we remove the unseen object in the testing set (*pitcher base*) when calculating MPCPE.

Under the symmetry model (denoted as "sym"), we re-

port the performance of our ArtiBoost in Tab. 3. Our *Clas* outperforms the recent Transformer-based method [23] when using ArtiBoost. We also evaluate our method on latest released dataset DexYCB and HO3D v3 in Tab. 4 and Tab. 5, respectively. All the results demonstrate the effectiveness of our method.

## 4.3. Ablation Study

To further discover how ArtiBoost works, we design two ablation studies. **A).** We compare the ArtiBoost with conventional grasp synthesis methods to show the efficacy of the CCV-space. **B).** We compare ArtiBoost with an offline training scheme to demonstrate the efficiency of our dynamic online re-weighting.

**A). Conventional Grasp Synthesis.** Simulated hand-object poses in GraspIt [48] are not necessarily correct or diverse (see Sec. 3.1 and Fig. 3 of GanHand [10]). Therefore, Corona *et al.* [10] manually annotated the MANO hand grasping YCB objects [4] and released a grasping pose dataset: YCBAfford. We find YCBAfford is an ideal contrast of our synthetic grasps in CCV-space. In this study, we compare the performance of the same model trained on two different data compositions. One is **a.1)** HO3D plus YCBAfford, and the other is **a.2)** HO3D plus our synthetic poses in CCV-space. To ensure fair and instructive comparison, we use the same amounts of HO poses from YCBAfford and CCV-space, set up an identical rendering pipeline, and turn off the re-weighting. During training, the synthetic HO poses will be randomly sampled and rendered, and then

| Training set composition | MPJPE | MPCPE |
|---|---|---|
| HO3D | 3.06 | 7.24 |
| **a.1)** HO3D $\bigoplus$ YCBAfford | 3.01 | 6.89 |
| **a.2)** HO3D $\bigoplus$ CCV-space | 2.71 | 5.49 |
| HO3D + Arti (full version) | **2.64** | **5.16** |

Table 6. Ablation on grasp poses synthesis. **a.1)** *v.s* **a.2)** shows a same network model (*Clas*) trained on HO3D plus synthetic pose from Conventional *v.s* from CCV-space. (*cm*)

| Method on % of source data | MPJPE | MPCPE |
|---|---|---|
| Our *Reg* (10%) | 3.81 | 8.77 |
| Our *Reg* (100%) | 3.53 | 7.38 |
| Our *Reg* (10%) **+ Arti** | **3.29** | **6.87** |
| Our *Clas* (10%) | 3.63 | 7.66 |
| Our *Clas* (100%) | 3.06 | 7.24 |
| Our *Clas* (10%) **+ Arti** | **3.05** | **6.02** |

Table 7. Performance of models trained on 10% of **HO3D** source data. (*cm*)

| Dataset | Method | MPJPE | MPCPE | CS-J |
|---|---|---|---|---|
| HO3Dv1 | [24] | 5.75 | 9.61 | 6.24 |
| HO3Dv1 | [24] + **Arti** | **3.67** | **3.24** | **3.57** |
| HO3D | [24] | 3.69 | 12.38 | 5.52 |
| HO3D | [24] + **Arti** | **3.39** | **8.31** | **4.90** |

Table 8. Results of porting ArtiBoost to the model in Hasson *et al*. [24]. **CS-J**: the MPJPE in camera space; all in *cm*.
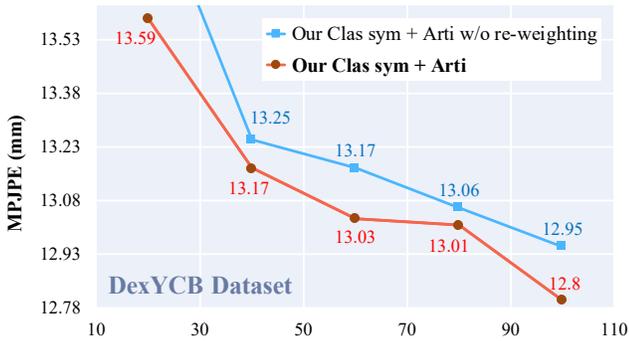


Figure 8. Ablation on offline training scheme.

the rendered images will be blended into the original HO3D training set. Tab. 6 shows the results of this study. We find that the model trained with poses in CCV-space outperforms the model trained with conventional synthetic poses, verifying an essential idea in our paper: diverse pose variants facilitate pose estimation.

**B). Offline Training Scheme.** To simulate the offline training scheme, we fix all the weights in the sampling weight map and randomly choose triplets from the CCV-space. We inspect two experiments throughout the entire training process, one of which uses the online sample re-weighting strategy, and the other follows the offline scheme. Both experiments use the *Clas* and symmetry model and are trained on the DexYCB dataset. We report the interim results on DexYCB testing set at certain intervals. As shown in Fig. 8, online sample re-weighting helps the models to converge fast and achieve a higher score.

### 4.4. Applications

We explore the potential application of ArtiBoost and design two studies. **A).** As real-world data labeling is inefficient and costly, we try to use ArtiBoost to help neural models on training with less amount of real-world labeled data. **B).** As ArtiBoost is model-agnostic, we show that it can be ported to other HOPE learning frameworks and boost their performance.

**A). Training on Less Real-world Labeled Data.** This study trains the neural models using only a small portion of the HO3D training data. Supposing the original amount of data in the HO3D training set is $N$, we set up three different amounts of training set: (1) 10% $N$ of the original set, (2) 100% $N$ of the original set, and (3) 10% $N$ of the original

set plus 100% $N$ of ArtiBoost synthetic poses. As shown in Tab. 7, we find that neural networks trained with setting (3) can outperform the network trained with 100% real-world data.

**B). Porting ArtiBoost to other HOPE Model.** In Tab. 8, we provide the results of porting ArtiBoost to another HOPE framework proposed by Hasson *et al*. [24] which directly regressed the hand-object poses and focal-normalized camera-space translations. The source training dataset: **HO3Dv1** [20] used in [24] is an early version of HO3D. We reproduced the results in [24] by training their network on the predefined **v1** set only. We report both the MPJPE in camera space and wrist-relative system. We show in Tab. 8 that porting ArtiBoost into a camera-space HOPE model significantly improves all metrics.

## 5. Discussion

**Limitation.** However, we do not explicitly mitigate the domain gap between the synthetic and real data, as we find that the dominant improvement to the HOPE task is brought from the images of more diverse pose variants, rather than images with a more realistic appearance. Besides, as the renderer in ArtiBoost is not differentiable, current ArtiBoost only supports exploration in a predefined lookup table (*e.g.* CCV-space). In future work, we will investigate a powerful generative and contrastive model seeking common features shared by both real and synthetic images.

**Conclusion.** In this work, we propose a novel online data enrichment method ArtiBoost, which enhances the learning framework of articulated pose estimation by exploration and synthesis. Our proposed ArtiBoost can be integrated into any learning framework, and in this work, we show its efficacy on the challenging task of hand-object pose estimation. Even with a simple baseline, our method can boost it to outperform the previous SOTA on the popular datasets. Besides, the proposed CCV-space also opens the door towards the generic articulated pose estimation, which we leave as future work.

# References

[1] Adnane Boukhayma, Rodrigo de Bem, and Philip HS Torr. 3d hand shape and pose from images in the wild. In *CVPR*, 2019. 2

[2] Samarth Brahmbhatt, Ankur Handa, James Hays, and Dieter Fox. ContactGrasp: Functional Multi-finger Grasp Synthesis from Contact. In *IROS*, 2019. 1, 2

[3] 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 *ECCV*, 2020. 1

[4] Berk Calli, Arjun Singh, Aaron Walsman, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M. Dollar. The ycb object and model set: Towards common benchmarks for manipulation research. In *ICAR*, 2015. 7

[5] Zhe Cao, Ilija Radosavovic, Angjoo Kanazawa, and Jitendra Malik. Reconstructing hand-object interactions in the wild. In *ICCV*, 2021. 1, 2

[6] Yu-Wei Chao, Wei Yang, Yu Xiang, Pavlo Molchanov, Ankur Handa, Jonathan Tremblay, Yashraj S. Narang, Karl Van Wyk, Umar Iqbal, Stan Birchfield, Jan Kautz, and Dieter Fox. DexYCB: A benchmark for capturing hand grasping of objects. In *CVPR*, 2021. 1, 2, 6, 12

[7] Wenzheng Chen, Huan Wang, Yangyan Li, Hao Su, Zhenhua Wang, Changhe Tu, Dani Lischinski, Daniel Cohen-Or, and Baoquan Chen. Synthesizing training images for boosting human 3d pose estimation. In *3DV*, 2016. 2

[8] Xingyu Chen, Yufeng Liu, Chongyang Ma, Jianlong Chang, Huayan Wang, Tian Chen, Xiaoyan Guo, Pengfei Wan, and Wen Zheng. Camera-space hand mesh recovery via semantic aggregationand adaptive 2d-1d registration. In *CVPR*, 2021. 1, 2

[9] Xingyu Chen, Yufeng Liu, Dong Yajiao, Xiong Zhang, Chongyang Ma, Yanmin Xiong, Yuan Zhang, and Xiaoyan Guo. MobRecon: Mobile-friendly hand mesh reconstruction from monocular image. In *arXiv preprint arXiv:2112.02753*, 2021. 1, 2

[10] Enric Corona, Albert Pumarola, Guillem Alenya, Francesc Moreno-Noguer, and Grégory Rogez. Ganhand: Predicting human grasp affordances in multi-object scenes. In *CVPR*, 2020. 1, 2, 7

[11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 12

[12] Bardia Doosti, Shujon Naha, Majid Mirbagheri, and David J Crandall. HOPE-Net: A graph-based model for hand-object pose estimation. In *CVPR*, 2020. 1

[13] Carlo Ferrari and John F Canny. Planning optimal grasps. In *ICRA*, 1992. 2

[14] Guillermo Garcia-Hernando, Shanxin Yuan, Seungryul Baek, and Tae-Kyun Kim. First-person hand action benchmark with rgb-d videos and 3d hand pose annotations. In *CVPR*, 2018. 1, 6

[15] Liuhao Ge, Zhou Ren, Yuncheng Li, Zehao Xue, Yingying Wang, Jianfei Cai, and Junsong Yuan. 3d hand shape and pose estimation from a single rgb image. In *CVPR*, 2019. 1, 2

[16] Weifeng Ge, Weilin Huang, Dengke Dong, and Matthew R. Scott. Deep metric learning with hierarchical triplet loss. In *ECCV*, 2018. 3

[17] Kehong Gong, Jianfeng Zhang, and Jiashi Feng. Poseaug: A differentiable pose augmentation framework for 3d human pose estimation. In *CVPR*, 2021. 3

[18] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The "something something" video database for learning and evaluating visual common sense. In *ICCV*, 2017. 1

[19] Patrick Grady, Chengcheng Tang, Christopher D Twigg, Minh Vo, Samarth Brahmbhatt, and Charles C Kemp. Contactopt: Optimizing contact to improve grasps. In *CVPR*, 2021. 1, 2

[20] Shreyas Hampali, Markus Oberweger, Mahdi Rad, and V. Lepetit. HO-3D: A multi-user, multi-object dataset for joint 3d hand-object pose estimation. In *arXiv preprint arXiv:1907.01481*, 2019. 8

[21] Shreyas Hampali, Mahdi Rad, Markus Oberweger, and Vincent Lepetit. Honnotate: A method for 3d annotation of hand and object poses. In *CVPR*, 2020. 1, 2, 6

[22] Shreyas Hampali, Sayan Deb Sarkar, and Vincent Lepetit. HO-3D_v3: Improving the accuracy of hand-object annotations of the ho-3d dataset. In *arXiv preprint arXiv:2107.00887*, 2021. 2, 6

[23] Shreyas Hampali, Sayan Deb Sarkar, Mahdi Rad, and Vincent Lepetit. HandsFormer: Keypoint transformer for monocular 3d pose estimation ofhands and object in interaction. In *arXiv preprint arXiv:2104.14639*, 2021. 6, 7, 12

[24] 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 *CVPR*, 2020. 1, 2, 7, 8

[25] Yana Hasson, Gül Varol, Ivan Laptev, and Cordelia Schmid. Towards unconstrained joint hand-object reconstruction from rgb videos. In *3DV*, 2021. 1, 2

[26] Yana Hasson, Gul Varol, Dimitrios Tzionas, Igor Kalevatykh, Michael J Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *CVPR*, 2019. 1, 2, 5

[27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6

[28] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017. 3

[29] Tomáš Hodaň, Martin Sundermeyer, Bertram Drost, Yann Labbé, Eric Brachmann, Frank Michel, Carsten Rother, and Jiří Matas. Bop challenge 2020 on 6d object localization. In *ECCV Workshops*, 2020. 7

[30] Lin Huang, Jianchao Tan, Jingjing Meng, Ji Liu, and Junsong Yuan. HOT-Net: Non-autoregressive transformer for 3d hand-object pose estimation. In *ACMMM*, 2020. 1

[31] Hanwen Jiang, Shaowei Liu, Jiashun Wang, and Xiaolong Wang. Hand-object contact consistency reasoning for human grasps generation. In *ICCV*, 2021. 3

[32] Korrawe Karunratanakul, Jinlong Yang, Yan Zhang, Michael J Black, Krikamol Muandet, and Siyu Tang. Grasping field: Learning implicit representations for human grasps. In *3DV*, 2020. 1, 3

[33] Wadim Kehl, Fabian Manhardt, Federico Tombari, Slobodan Ilic, and Nassir Navab. SSD-6d: Making rgb-based 3d detection and 6d pose estimation great again. In *ICCV*, 2017. 2

[34] Mia Kokic, Danica Kragic, and Jeannette Bohg. Learning task-oriented grasping from human activity datasets. *IEEE Robotics and Automation Letters*, 5(2), 2020. 1, 2

[35] Felix Kuhnke and Jörn Ostermann. Deep head pose estimation using synthetic images and partial adversarial domain adaption for continuous label spaces. *ICCV*, 2019. 3

[36] Taein Kwon, Bugra Tekin, Jan Stuhmer, Federica Bogo, and Marc Pollefeys. H2o: Two hands manipulating objects for first person interaction recognition. In *ICCV*, 2021. 1

[37] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang, and Cewu Lu. HybrIK: A hybrid analytical-neural inverse kinematics solution for 3d human pose and shape estimation. In *CVPR*, 2021. 1, 2

[38] Xiaolong Li, He Wang, Li Yi, Leonidas J Guibas, A Lynn Abbott, and Shuran Song. Category-level articulated object pose estimation. In *CVPR*, 2020. 1, 2

[39] John Lin, Ying Wu, and Thomas S Huang. Modeling the constraints of human hand motion. In *Proceedings workshop on human motion*, 2000. 3

[40] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 3

[41] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 5

[42] Liu Liu, Han Xue, Wenqiang Xu, Haoyuan Fu, and Cewu Lu. Towards real-world category-level articulation pose estimation. *IEEE Transactions on Image Processing*, 2022. 1, 2

[43] Shaowei Liu, Hanwen Jiang, Jiarui Xu, Sifei Liu, and Xiaolong Wang. Semi-supervised 3d hand-object poses estimation with interactions in time. In *CVPR*, 2021. 1, 2, 7

[44] Jun Lv, Wenqiang Xu, Lixin Yang, Sucheng Qian, Chongzhao Mao, and Cewu Lu. HandTailor: Towards high-precision monocular 3d hand recovery. In *BMVC*, 2021. 7

[45] George Marsaglia et al. Choosing a point from the surface of a sphere. *The Annals of Mathematical Statistics*, 1972. 5

[46] Joanna Materzynska, Tete Xiao, Roei Herzig, Huijuan Xu, Xiaolong Wang, and Trevor Darrell. Something-else: Compositional action recognition with spatial-temporal interaction networks. In *CVPR*, 2020. 1

[47] Matthew Matl. PyRender. https://github.com/mmatl/pyrender, 2019. 5

[48] A.T. Miller and P.K. Allen. Graspit! a versatile simulator for robotic grasping. *IEEE Robotics Automation Magazine*, 2004. 2, 7

[49] Franziska Mueller, Florian Bernard, Oleksandr Sotnychenko, Dushyant Mehta, Srinath Sridhar, Dan Casas, and Christian Theobalt. Ganerated hands for real-time 3d hand tracking from monocular rgb. In *CVPR*, 2018. 5

[50] Neng Qian, Jiayi Wang, Franziska Mueller, Florian Bernard, Vladislav Golyanik, and Christian Theobalt. Html: A parametric hand texture model for 3d hand reconstruction and personalization. In *ECCV*, 2020. 5

[51] Gregory Rogez, James S. Supancic, III, and Deva Ramanan. Understanding everyday hands in action from rgb-d images. In *ICCV*, 2015. 1, 2

[52] Javier Romero, Hedvig Kjellström, and Danica Kragic. Hands in action: real-time 3d reconstruction of hands in interaction with objects. In *ICRA*, 2010. 1, 2

[53] Javier Romero, Dimitrios Tzionas, and Michael J Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics*, 2017. 2, 3

[54] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015. 3

[55] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *CVPR*, 2016. 3

[56] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *NIPS*, 2016. 3

[57] Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. Integral human pose regression. In *ECCV*, 2018. 2

[58] Martin Sundermeyer, Zoltan-Csaba Marton, Maximilian Durner, Manuel Brucker, and Rudolph Triebel. Implicit 3d orientation learning for 6d object detection from rgb images. In *ECCV*, 2018. 2

[59] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. GRAB: A dataset of whole-body human grasping of objects. In *ECCV*, 2020. 3, 4, 5

[60] Xiao Tang, Tianyu Wang, and Chi-Wing Fu. Towards accurate alignment in real-time 3d hand-mesh reconstruction. In *ICCV*, 2021. 2

[61] Bugra Tekin, Federica Bogo, and Marc Pollefeys. H+O: Unified egocentric recognition of 3d hand-object poses and interactions. In *CVPR*, 2019. 1, 6

[62] Bugra Tekin, Sudipta N Sinha, and Pascal Fua. Real-time seamless single shot 6d object pose prediction. In *CVPR*, 2018. 2

[63] Gül Varol, Ivan Laptev, Cordelia Schmid, and Andrew Zisserman. Synthetic humans for action recognition from unseen viewpoints. *International Journal of Computer Vision*, 2021. 1, 2

[64] Chengde Wan, Thomas Probst, Luc Van Gool, and Angela Yao. Self-supervised 3d hand pose estimation through training by fitting. In *CVPR*, 2019. 2

[65] Haonan Yan, Jiaqi Chen, Xujie Zhang, Shengkai Zhang, Nianhong Jiao, Xiaodan Liang, and Tianxiang Zheng. Ultrapose: Synthesizing dense pose with 1 billion points by human-body decoupling 3d model. In *ICCV*, 2021. 1, 2

[66] Linlin Yang, Shicheng Chen, and Angela Yao. Semihand: Semi-supervised hand pose estimation with consistency. In *ICCV*, 2021. 2

[67] Lixin Yang, Jiasen Li, Wenqiang Xu, Yiqun Diao, and Cewu Lu. BiHand: Recovering hand mesh with multi-stage bisected hourglass networks. In *BMVC*, 2020. 2

[68] Linlin Yang and Angela Yao. Disentangling latent hands for image synthesis and pose estimation. In *CVPR*, 2019. 2

[69] Lixin Yang, Xinyu Zhan, Kailin Li, Wenqiang Xu, Jiefeng Li, and Cewu Lu. CPF: Learning a contact potential field to model the hand-object interaction. In *ICCV*, 2021. 1, 2, 4

[70] Yi Zhou, Connelly Barnes, Lu Jingwan, Yang Jimei, and Li Hao. On the continuity of rotation representations in neural networks. In *CVPR*, 2019. 12

[71] Yuxiao Zhou, Marc Habermann, Weipeng Xu, Ikhsanul Habibie, Christian Theobalt, and Feng Xu. Monocular real-time hand shape and motion capture using multi-modal data. In *CVPR*, 2020. 2

[72] Christian Zimmermann and Thomas Brox. Learning to estimate 3d hand pose from single rgb images. In *ICCV*, 2017. 1, 2, 5

# Appendices

## A. The Training Details

The backbones of classification-based (*Clas*) and regression-based (*Reg*) baseline networks are initialized with ImageNet [11] pretrained model. In *Clas*, the output resolution of 3D-heatmaps is $28 \times 28 \times 28$. The MLP branch that predicts object rotation adopts three fully-connected layers with 512, 256 and 128 neurons for each, and a final layer of 6 neurons that predict the continuity representation [70] of object rotation: $\mathbf{r}_o \in \mathfrak{so}(3)$. We train the network 100 epochs with Adam optimizer and learning rate of $5 \times 10^{-5}$. The training batch size across all the following experiments is 64 per GPU and 2 GPUs in total. The framework is implemented in PyTorch. All the object models and textures are provided by the original dataset. For all the training batches, the blended rate of original real-world data and ArtiBoost synthetic data is approximately $1:1$. We empirically find that this real-synthetic blended rate achieves the best performance.

## B. Objects' Symmetry Axes

In the hand-object interaction dataset, it is far more challenging to predict the pose of an object than in the dataset that only contains objects, since the objects are often severely occluded by the hand. Therefore, we relax the restrictions of the objects' symmetry axes following the practices in [6, 23]. Supposing the set $\mathcal{S}$ contains all the valid rotation matrices based on the object's predefined symmetry axes, we calculate $\mathcal{S}$ with the following step:

| Objects | Axes: $\mathbf{n}$ | Angle: $\theta$ |
|---|---|---|
| 002_master_chef_can | x, y, z | $180°, 180°, \infty$ |
| 003_cracker_box | x, y, z | $180°, 180°, 180°$ |
| 004_sugar_box | x, y, z | $180°, 180°, 180°$ |
| 005_tomato_soup_can | x, y, z | $180°, 180°, \infty$ |
| 006_mustard_bottle | z | $180°$ |
| 007_tuna_fish_can | x, y, z | $180°, 180°, \infty$ |
| 008_pudding_box | x, y, z | $180°, 180°, 180°$ |
| 009_gelatin_box | x, y, z | $180°, 180°, 180°$ |
| 010_potted_meat_can | x, y, z | $180°, 180°, 180°$ |
| 024_bowl | z | $\infty$ |
| 036_wood_block | x, y, z | $180°, 180°, 90°$ |
| 037_scissors | z | $180°$ |
| 040_large_marker | x, y, z | $180°, \infty, 180°$ |
| 052_extra_large_clamp | x | $180°$ |
| 061_foam_brick | x, y, z | $180°, 90°, 180°$ |

Table 9. **YCB objects' axes of symmetry**. $\infty$ indicates the object is revolutionary by the axis.

1) Firstly, as shown in Fig 9, we align the object to its principal axis of inertia.
2) Secondly, we define the axis $\mathbf{n}$ and angle $\theta$ of symmetry in Tab 9 under the aligned coordinate system, where the object's geometry does not change when rotate this object by an angle of $\theta$ around $\mathbf{n}$. Here we get the predefined rotation matrix $\mathbf{R}_{def} = \exp(\theta\mathbf{n})$.
3) To get a more accurate rotation matrix $\mathbf{R}$, we use the Iterative Closest Point (ICP) algorithm to fit a $\Delta\mathbf{R}$. The ICP minimizes the difference between $\Delta\mathbf{R} * \mathbf{R}_{def} * \mathbf{V}_o$ and $\mathbf{V}_o$, where $\mathbf{V}_o$ is the point clouds on object surface. Finally, we have $\mathbf{R} = \Delta\mathbf{R} * \mathbf{R}_{def}, \mathbf{R} \in \mathcal{S}$.



Figure 9. **YCB objects' principal axis of inertia.** The x, y and z axis are colored in red, green and blue, respectively.

## C. Additional Results

We demonstrate 20 YCB objects' MSSD on DexYCB in Tab. 10. With ArtiBoost, our network can predict a more accurate pose for almost every object. More qualitative results on HO3D and DexYCB testing set are shown in Fig. 10.

| Objects | Our *Clas* sym | Our *Clas* sym **+ Arti** | Objects | Our *Clas* sym | Our *Clas* sym **+ Arti** |
|---|---|---|---|---|---|
| 002_master_chef_can | 27.62 | **25.59** | 003_cracker_box | 63.68 | **46.13** |
| 004_sugar_box | 48.42 | **39.20** | 005_tomato_soup_can | 33.31 | **31.90** |
| 006_mustard_bottle | 35.16 | **32.01** | 007_tuna_fish_can | 24.54 | **23.81** |
| 008_pudding_box | 39.92 | **35.04** | 009_gelatin_box | 45.99 | **37.81** |
| 010_potted_meat_can | 41.44 | **36.47** | 011_banana | 98.69 | **79.87** |
| 019_pitcher_base | 105.66 | **84.82** | 021_bleach_cleanser | 91.66 | **72.31** |
| 024_bowl | **31.74** | 32.37 | 025_mug | 65.46 | **54.28** |
| 035_power_drill | 74.95 | **52.70** | 036_wood_block | 51.24 | **50.69** |
| 037_scissors | 88.10 | **66.52** | 040_large_marker | 30.76 | **29.33** |
| 052_extra_large_clamp | 78.87 | **55.87** | 061_foam_brick | 34.23 | **31.53** |

Table 10. Full MSSD results ($mm$) on **DexYCB** testing set.



Figure 10. (**Best view in color**) More qualitative results on **HO3D** ($1^{st} \sim 3^{rd}$ rows) and **DexYCB** ($4^{th} \sim 8^{th}$ rows) datasets.