# D3D-HOI: Dynamic 3D Human-Object Interactions from Videos

Xiang Xu[1*]   Hanbyul Joo[2]   Greg Mori[1]   Manolis Savva[1]
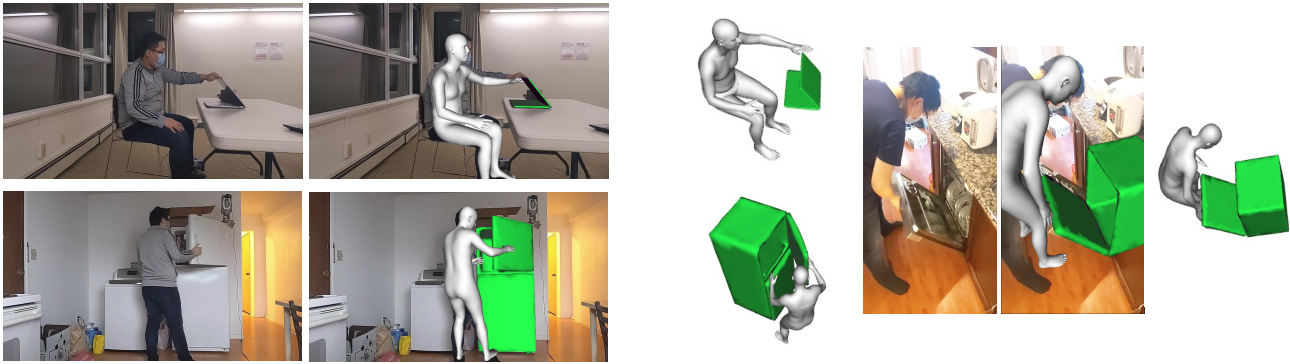
[1]Simon Fraser University   [2]Facebook AI Research

Figure 1: We address the problem of recovering dynamic 3D human-object interactions given a video input. Our focus is on reconstructing the articulating 3D object during manipulation. To enable research in this direction, we collect a dataset of interaction videos with common objects (e.g., laptop, fridge, dishwasher). We annotate the 3D object pose, shape, articulation state, and estimate the 3D mesh of the manipulator using the approach from Joo et al. [10]. Above are rendered frames from the ground truth annotations.

## Abstract

*We introduce D3D-HOI: a dataset of monocular videos with ground truth annotations of 3D object pose, shape and part motion during human-object interactions. Our dataset consists of several common articulated objects captured from diverse real-world scenes and camera viewpoints. Each manipulated object (e.g., microwave oven) is represented with a matching 3D parametric model. This data allows us to evaluate the reconstruction quality of articulated objects and establish a benchmark for this challenging task. In particular, we leverage the estimated 3D human pose for more accurate inference of the object spatial layout and dynamics. We evaluate this approach on our dataset, demonstrating that human-object relations can significantly reduce the ambiguity of articulated object reconstructions from challenging real-world videos. Code and dataset are available at https://github.com/facebookresearch/d3d-hoi.*

## 1. Introduction

3D reconstruction has gained prominence in recent years. The goal is to infer the full 3D shape, pose and layout from only partial 2D information. This is challenging due to the inherent ambiguities caused by mapping back from 2D to 3D. Most works in this area focus on separately reconstructing the human mesh [1, 10, 11, 15, 16, 27] or the individual 3D object [4, 18, 22, 31]. Recent works [5, 29, 32, 35] do explore human-object relations for more accurate 3D reconstruction, but they only deals with static objects with no part motion.

Real world human-object interactions (HOIs) often involves dynamic articulated objects (e.g., opening and closing a fridge). Traditional approaches [7, 18, 33] for articulated object reconstruction requires accurate RGB-D data and do not leverage the human-object relations. In response to the lack of 3D HOI data, we first create the Dynamic 3D HOI (D3D-HOI) dataset. We then provide an optimization-based method that utilize human-object relations to reconstruct the articulated objects from only RGB video.

D3D-HOI is a real world HOI video dataset with 3D

---

annotations of object location, orientation, size, matching CAD model, and part motion state at every video frame. Our annotations are based on 3D mesh models from Xiang et al. [33] which we use to represent objects with moving parts that articulate along a specific motion axis and origin. We focus on translation (prismatic joint) and rotation (revolute joint) motions. For simplicity, we only consider one joint motion per object part. Note that an object can have multiple parts and each part will have its own unique joint motion. Figure 1 illustrates rendered examples of the ground truth annotations. D3D-HOI allows us to evaluate reconstruction methods along several axes: shape, pose and part motion of real-world articulated objects during human-object interaction.

Using the D3D-HOI dataset, we explore the role of human-object relations and propose an optimization method that leverages the estimated human pose and dynamics to better reconstruct the interacted object. Our insight is that treating both the human and object as dynamic entities allow us to constrain the pose and motion of each through orientation and contact terms. To the best of our knowledge, we are the first to reconstruct articulated object from real-world HOI video.

In summary, we make the following two contributions:

- We collect a video dataset of human interactions with articulated objects and provide the ground truth annotations of object 3D pose, shape and part motion at every time frame.

- We present an optimization method based on human-object relations for the task of articulated object reconstruction from HOI video. We also evaluate the 3D reconstruction accuracy on our new dataset.

## 2. Related Work

### 2.1. 3D human pose and motion estimation

Recently there has been a lot of interest for estimating human pose and motion from 2D image or video. Fitting-based methods assume a parametric body model such as SMPL [20] or SMPL-X [27] and use optimization algorithms like SMPLify [1] to fit the parameters. On the other hand, regression-based methods [11, 25, 26] rely on deep neural networks and large amounts of training data [8, 23] to directly predict the 3D pose of the human. Hybrid methods [10, 16] also demonstrate improved results. Compared to image-based methods, video-based methods [15, 21] can utilize temporal information and better estimate human pose and motion over entire video sequences.

### 2.2. 3D object pose and motion estimation

Recovering object shape and pose from 2D images is a very active field. Recent methods [4, 17, 24, 31] can achieve impressive results on indoor environments. These methods

are designed to work on non-articulated objects where object state is fixed and there is no part motion.

Methods that do predict part motion parameters usually come from the field of embodied ai, where the goal is to teach robots how to open or close doors. Synthetic datasets [30, 33] are used to train regression-based models to directly predict part motion parameters. However, methods train on these dataset require accurate depth information [7, 9, 18, 34]. Thus it is not straightforward for one to directly apply them to real world scenes due to noisy depth estimation and domain gap differences. Also none of these methods utilize human information during HOI.

### 2.3. 3D human-object interaction

Some recent works do examine 3D human-object relations during HOI and use them to improve performance on various tasks. PiGraphs [29] collects a dataset with 3D annotation of human and interacted objects (i.e. chair, desk, whiteboard). Hassan et al. [5] introduced the PROX dataset which consists of 3D human mesh and static 3D scene. They use static 3D scene structure to improve human pose estimation from monocular images. Follow-up works [6, 36, 37] also use this dataset to generate human mesh and insert it to an environment based on 3D scene context. Human-object relations are also important for 3D reconstruction. In [3], the authors build a graph based on human-object relations and propose a holistic scene parsing and reconstruction pipeline. Zhang [35] jointly infers spatial layout and shapes of humans and objects in a consistent 3D scene. Karunratanakul [12] defines a grasping field to jointly reconstructs the human hand and interacted object. For those works, only the human has motion and the object remains static. This is different from our dynamic human-object interaction setting where both object and human can have motions.

## 3. Dynamic 3D HOI Reconstruction

Given a video from a fixed camera viewpoint where a human dynamically interacts with an articulated object, we aim at estimating the 3D human motion, 3D movement of the articulated object, and their 3D arrangement. In particular, we use articulated parametric models to represent humans and objects respectively. This takes into account the dynamic interactions between them while ignoring detailed shape variations. Specifically, we use the SMPL model [20] to represent the 3D state of a human at time $t$ by $\Theta_h^t = (\boldsymbol{x}_h^t, \boldsymbol{\theta}_h^t, \boldsymbol{\beta}_h)$, where $\boldsymbol{x}_h^t \in \mathbb{R}^3$ is 3D translation, $\boldsymbol{\theta}_h^t \in \mathbb{R}^{24 \times 3}$ controls the rotations of 24 body joints with respect to their parent joints, and the shape $\beta_h \in \mathbb{R}^{10}$ controls the shape of the body. Similarly, we represent each 3D object at time $t$ via a parametric model $\Theta_o^t = (\boldsymbol{x}_o, \boldsymbol{\phi}_o, \boldsymbol{\theta}_o^t, \beta_o)$,

|              | refrigerator | storageFurniture | trashcan | washingmachine | microwave | dishwasher | laptop | oven |
|--------------|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| video        | 82  | 32  | 19  | 24  | 36  | 11  | 41  | 11  |
| cad model    | 7   | 4   | 2   | 3   | 2   | 2   | 2   | 2   |
| viewpoint    | 14  | 5   | 4   | 7   | 7   | 4   | 8   | 3   |
| scene        | 5   | 4   | 3   | 3   | 2   | 2   | 2   | 1   |

Table 1: Statistics of the Dynamic 3D HOI Dataset with category-level distributions for number of videos, 3D CAD models, camera viewpoints, and 3D scenes.

where $\boldsymbol{x}_o \in \mathbb{R}^3$ is the global translation and $\boldsymbol{\phi}_o \in \mathbb{R}^{3\times3}$ is the global orientation of the object base. In particular, $\boldsymbol{\theta}_o^t \in \mathbb{R}^n$ represents the 3D part motion value of the dynamic part. This is either rotation degree for revolute joint or translation offset for prismatic joint. The object shape parameter $\boldsymbol{\beta}_o \in \mathbb{R}^3$ represents the width, length, and height sizes of the object. Thus, our target task is to estimate $\{\Theta_h^t\}_{t=[1,n]}$ and $\{\Theta_o^t\}_{t=[1,n]}$ from a monocular video, $\{I_t\}_{t=[1,n]}$. Note that $\Theta_h^t$ and $\Theta_o^t$ are defined in the same 3D world space. Examples are shown in Fig 1. Since there has been a lot of works for 3D human reconstruction using SMPL, we use an existing method [10] to extract the 3D human shape and pose. The focus of this work is thus on the estimating $\{\Theta_o^t\}_{t=[1,n]}$.

## 4. Dynamic 3D HOI Dataset

A major obstacle in pursuing the study of dynamic 3D HOI reconstruction task is the lack of available datasets on which algorithms are quantitatively compared and evaluated. To tackle this challenge, we create the Dynamic 3D HOI Dataset (D3D-HOI), providing videos of dynamically interacting individuals[1], with manually annotated ground-truth 3D object parameters, $\{\Theta_o^t\}_{t=[1,n]}$. In this section, we describe the statistics and summary of our dataset (subsection 4.1), the data processing procedures to build 3D articulated models (subsection 4.2), and the data collection pipeline for video capture and annotations (subsection 4.3),

### 4.1. Data Statistics

D3D-HOI contains a total of 256 videos. After subsampling videos into 3 FPS, our dataset has 6286 image frames in total. Each frame is annotated with 3D object dimension, matching CAD model, location, orientation, and part motion. Data acquisition involved 5 volunteers interacting with objects in 22 different scenes. The articulated objects come from 8 SAPIEN categories and were represented by a total of 24 CAD models. The camera is placed at 52 different viewpoints from the object. Table 1 summarizes the overall dataset statistics. From the summary we see that the col-

lected dataset contains a variety of articulated objects from small (e.g., laptop, microwave) to large (e.g., refrigerator, dishwasher). The distribution over viewpoints, scenes, and 3D models are also very diverse and covers a wide range of human-object interaction scenarios. Example annotated frames from our dataset is shown in Figure 1.

### 4.2. Articulated Models

We consider object categories that are commonly observable in our daily life and have movable parts with which humans can interact: refrigerator, storage furniture (including drawer), microwave, dishwasher, refrigerator, trash can, washing machine, laptop, and oven. To build the parametric models for articulated objects, we leverage the SAPIEN PartNet-Mobility [33] dataset by applying additional post-processing to simplify mesh topology. The PartNet-Mobility dataset provides various 3D articulated mesh models with movable part definitions. Among the models in our target object categories, we choose the 3D CAD models (up to four per category) most similar to the objects captured in our input videos. We also rescale the 3D models to have the same size with the actual objects we captured in measurement units (height, width, and length in $cm$). We define a local object coordinate frame by positioning the origin at the center of mesh and by making the "frontal" side of the object to be aligned with the $+z$ direction in the local coordinate. The $+y$ direction is defined as the ground normal direction when object is normally sitting on the floor. Finally, we simplify mesh topology to contain roughly 2500 triangles per object. Figure 2 shows examples of processed PartNet-Mobility models and their available part motions.

### 4.3. Data Collection Pipeline

**Capture System** We use two cameras (Google Pixel 4 and AVerMedia live streamer 313) to simultaneously record RGB videos on the HOI scenes from two different view points. One camera (Pixel 4) records the scenes from a fixed third person view where the whole human body is visible, and another camera (AVerMedia) is placed on the object to precisely record the movement of object part with a measurement tools (e.g., protractor). See Figure 3. This multi-camera setup is helpful to manually annotate the parame-

---

[1]The videos have been captured by the authors, with approvals from the individuals appearing in the videos to release the data for research purposes.

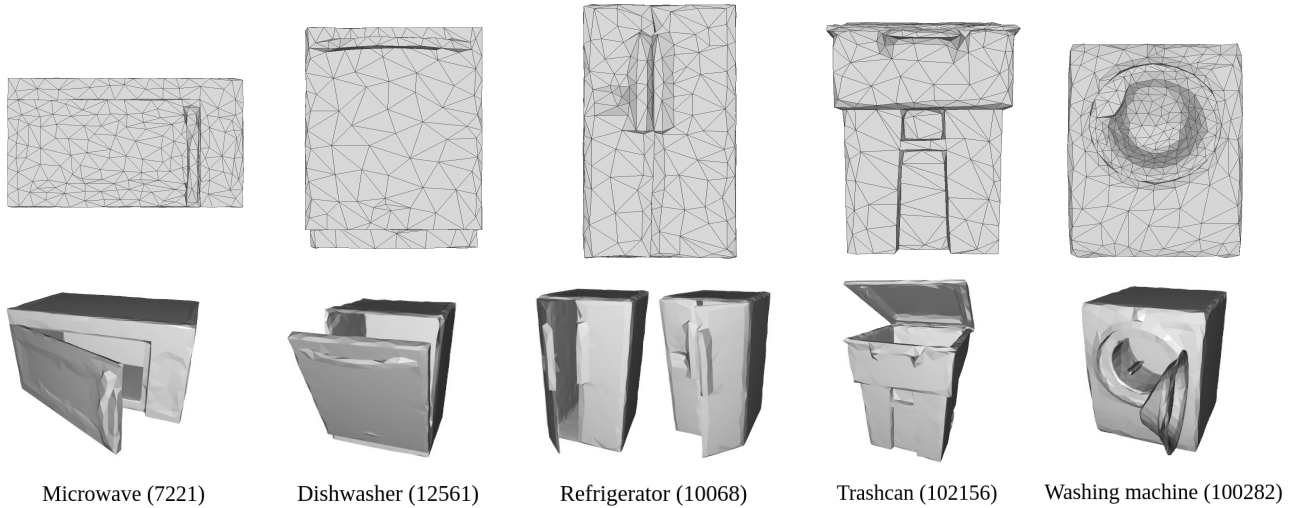Microwave (7221)  Dishwasher (12561)  Refrigerator (10068)  Trashcan (102156)  Washing machine (100282)

Figure 2: Examples of post-processed meshes from the PartNet-Mobility dataset [33]. Top row shows the simplified meshes with centering and front direction in $+z$ direction. Bottom row illustrates the available part motions for each object.



Figure 3: Left is an example of our equipment setup with dual cameras (one for recording, one for reading motion parameters). Right shows an example of the synchronized raw frame.

ters of movable parts (i.e. angles in degrees) at each video frame. All videos are recorded at 30 FPS.

**Calibration**  To estimate the intrinsic camera parameters, we perform camera calibration for each camera we used with a checkerboard pattern using OpenCV [2]. Since the front camera of Pixel 4 and the AVerMedia live streamer 313 have fixed focus, the calibration parameters are valid for all captured videos. However, for the rear Pixel 4 camera with auto-focus, our calibration parameters estimated afterwards may not be precise. To check this issue, we perform camera calibration multiple times by locking the focus with various target object distances similar to our capture setup (1.5 to 3.0 meters away from the camera). We compared the parameters, and confirmed that after locking auto-focus has minor intrinsic parameter differences. To this end, we use the averaged focal length estimated at different distance as
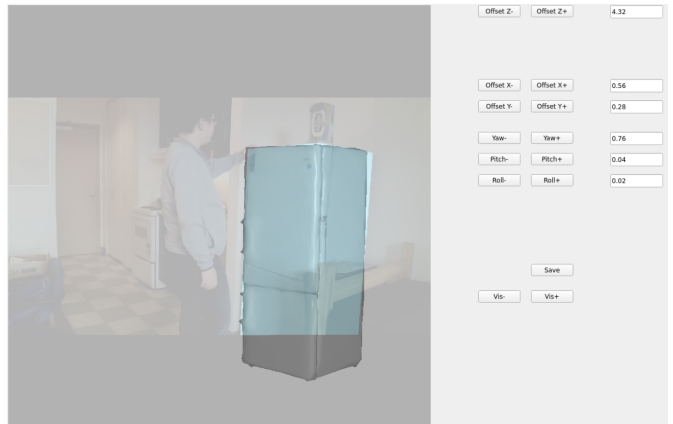


Figure 4: Interface for annotating the object location and orientation values.

the intrinsic parameters for this camera.

**Annotation Interface**  We build an interactive GUI interface for annotating the 3D orientation and location of the object. A screenshot of the interface is shown in Figure 4. In our annotation interface, users can see the input image frames together with the projection of 3D objects using the current object parameters. For the projection, the calibrated intrinsic parameters are used. Users can manually adjust the object parameters: translation parameter $x_o$ and orientation parameter $\phi_o$. The 3D object is re-rendered with the adjusted parameters in real-time so that users can confirm that the object projection is aligned to the input image. We save the final values as ground-truth annotation data. Note

that these translation and orientation parameters are fixed across frames, since the videos are recorded from a fixed camera view point and objects are static except the movable parts. The parameters for the movable part (e.g. doors) are separately annotated by observing the videos recorded with measurement tool (e.g. protractor) and annotating the angles at 3 FPS.

## 5. Methods

We first describe our design of the differentiable articulated object in 6. We then go over the objective function in Section 7. Section 7.1 provides detailed description of the object mask error. Section 7.2 shows how we fit human mesh to the same space as the object. Human-object interaction error is discussed in Section 7.3. Details of object parameter regularization is in Section 7.4. And finally we discuss the optimization pipeline in Section 8.

## 6. Differentiable Articulated Object

In this section, we describe how we differentiate from 3D object mesh to the articulated object parameters. Let us denote the processed object mesh $\mathbf{V} = \{\mathbf{v}_i\}$ where $\mathbf{v}_i \in \mathbb{R}^3$ is the $i$-th vertex location at the local object coordinate (local coordinate with origin at the center). The object mesh can be divided into movable part and non-movable part (base): $\mathbf{V} = \mathbf{V}_{\text{part}} \cup \mathbf{V}_{\text{base}}$. Given object parameters $\Theta_o^t = (\boldsymbol{x}_o, \boldsymbol{\phi}_o, \boldsymbol{\theta}_o^t, \beta_o)$ at time $t$, the object mesh is transformed as:

$$\hat{\mathbf{V}}(\Theta_o^t) = \{R_{\boldsymbol{\phi}_o}\left(\Psi_{\boldsymbol{\theta}_o^t}(\beta_o \mathbf{v}_i)\right) + \boldsymbol{x}_o\}, \quad (1)$$

where $\beta_o \in \mathbb{R}^3$ is scaling factor, $R_{\boldsymbol{\phi}_o}$ is the global rotation matrix, and $\boldsymbol{x}_o$ is the global translation. $\Psi_{\boldsymbol{\theta}_o^t}$ represents the part mobility matrix which is dependent on the selected parametric model and part motion value at time $t$. $\Psi_{\boldsymbol{\theta}_o^t}$ is the identity matrix for $\mathbf{v}_i \in \mathbf{V}_{base}$ since the object base is static with no motion.

For revolute motion, $\Psi_{\boldsymbol{\theta}_o^t}$ is a rotation matrix defined by the rotation axis of corresponding SAPIEN model together with the part rotation angle $\boldsymbol{\theta}_o^t$ at time $t$. Note that matrix $\Psi_{\boldsymbol{\theta}_o^t}$ is fully differentiable w.r.t value $\boldsymbol{\theta}_o^t$. For prismatic-joint, we follow a similar approach and define $\Psi_{\boldsymbol{\theta}_o^t}$ to be a transformation matrix correspond to the given translation axis and the offset value $\boldsymbol{\theta}_o^t$. In this way, gradient can back-prop from object mesh vertices back to the articulated object parameters.

## 7. Objective Function

To estimate the pose, part motion, and shape of the articulated object, we minimize an objective function that is the sum of four error terms: object mask term $E_{\text{m}}$, human fitting term $E_{\text{h}}$, human-object interaction term $E_{\text{hoi}}$, and a regularization term $E_{\text{r}}$. The overall objective $E(\Theta_o^t, \boldsymbol{x}_h^t, R_c) =$

$$\lambda_{\text{m}} E_{\text{om}}(\Theta_o^t) + \lambda_{\text{h}} E_{\text{h}}(\boldsymbol{x}_h^t; H_{\text{est}}^t) + \\ \lambda_{\text{hoi}} E_{\text{hoi}}(\Theta_o^t, R_c; \boldsymbol{x}_h^t, H_{\text{est}}^t) + \lambda_{\text{r}} E_{\text{r}}(\Theta_o^t). \quad (2)$$

Here $\lambda_{\text{m}}, \lambda_{\text{h}}, \lambda_{\text{hoi}}, \lambda_{\text{r}}$ are the hyper-parameters, object parameters $\Theta_o^t = (\boldsymbol{x}_o, \boldsymbol{\phi}_o, \boldsymbol{\theta}_o^t, \beta_o)$, $H_{\text{est}}^t$ is the estimated orthographic human mesh from [10], $\boldsymbol{x}_h^t$ are the translation parameters that transforms $H_{\text{est}}^t$ to the same 3D space as the object, and finally $R_c$ is the rotation matrix used for human-object contact curve matching. Next, we will describe those error terms in more details.

### 7.1. Object Mask Error

The object mask term penalizes inconsistency between the projected 3D object mask and the annotated object mask $M(t)$ estimated from PointRend [14]. It is based on a differentiable renderer [19] implemented in PyTorch3D [28]. Their method provides fully differentiable projection $P_{\text{soft}}$ from 3D mesh vertices to 2D image mask. Here, we set up $P_{\text{soft}}$ as perspective projection using our estimated focal length $K$. After applying global transformation and part motion, object meshes are passed through the differentiable render and we get the rendered 2D mask at every frame. The final object mask error is defined as:

$$E_{\text{m}}(\Theta_o^t) = \frac{1}{N} \sum_{t=1}^{N} ||P_{\text{soft}}\left(\hat{\mathbf{V}}(\Theta_o^t); K\right) - M(t)||^2, \quad (3)$$

where $\hat{\mathbf{V}}(\Theta_o^t)$ is defined as in Eq. 1.

### 7.2. Human Fitting Error

Human fitting error measures the difference between projected human mesh vertex and the estimated 2D vertex locations $p$ (also from EFT [10]) in the same frame. Minimizing the human fitting term $E_{\text{h}}$ will transform the estimated orthographic human body $H_{\text{est}}^t$ into the same 3D space as the object. We use a perspective projection $P_{\text{persp}}$ similar to $P_{\text{soft}}$ with the same focal length $K$ and scale the human mesh with a pre-determined scale $S_h$ so that the SMPL mesh height is equal to the height of the volunteer appearing in the video (170 cm for all videos). We then add 3D translation offset $\boldsymbol{x}_h^t$ to all the vertices at each time frame. The human fitting error is defined as $E_{\text{h}}(\boldsymbol{x}_h^t; H_{\text{est}}^t) =$

$$\frac{1}{N} \sum_{t=1}^{N} ||P_{\text{persp}}\left(S_h H_{\text{est}}^t + \boldsymbol{x}_h^t; K\right) - p(t)||^2. \quad (4)$$

Note that we allow per-frame 3D translation to be applied to the human mesh so that human can move around in the video.

### 7.3. Human-object Interaction Error

After inserting human into the same 3D world coordinate as the object, we now describe our human-object interaction term. This consists of two separate parts, orientation error and contact error. We will now describe each one in details.

**Orientation error** Our first human-object interaction term is the orientation error. During interaction, the human often faces towards the interacted object. Motivated by this observation, we estimate the human facing direction from SMPL and encourage object front direction to be the opposite. We define human facing direction $D_f(t)$ as the cross product of two vectors: SMPL left shoulder to SMPL right hip, and SMPL right shoulder to SMPL left hip. Object front direction is the vector pointing from the object center towards the direction where human interaction most frequently occurs. For our post-processed data, this direction is $+z$ in the local object coordinate before applying global rotation.

In addition to human facing direction, we also want to align human and object so that their inferred ground normal directions are parallel. We define ground normal direction $D_g(t)$ from human as the cross product of two vectors: left SMPL feet ankle to right SMPL feet toe, and right SMPL feet toe to left SMPL feet ankle. The object ground normal direction is always pointing towards $+y$ direction in the local object coordinate. The overall human orientation error is defined as $E_{\text{orientation}}(\Theta_o^t; \boldsymbol{x}_h^t, H_{\text{est}}^t) =$

$$\sum_{t=t_{\text{start}}}^{t_{\text{end}}} (cos(D_f(t), R_{\boldsymbol{\phi}_o}\vec{-z}) + cos(D_g(t), R_{\boldsymbol{\phi}_o}\vec{y})). \quad (5)$$

where $R_{\boldsymbol{\phi}_o}$ is the object global rotation matrix, $t_{\text{start}}, t_{\text{end}}$ are the start and end time of the interaction where the object part is moving.

**Contact error** Our second human-object interaction term is the contact error. We make the assumption that during human-object interaction, the object motions tends to follow the human hand. This allows us to use the 3D hand location from the human mesh to constrain the 3D location of moving part.

One issue that arises is that it is difficult to find the exact contact vertex on the object when human is touching it. During optimization, object orientation and location in general do not closely match the ground-truth. Matching the 3D human-object contact location introduce noise and might even lead to decreased performance. To solve this issue, we introduce a pose-invariant contact term that can tolerate such noise by only matching the 3D shape of the human-object contact curves. We select key vertex locations on the 3D CAD models which are most representative

of the shape of the moving part during interaction. During optimization, we match against both the left and right hands with the selected object contact vertex and report the final result with the lowest optimization cost.

Assume our selected object contact vertex is $v_c$. We then have the same vertex in 3D camera coordinate (world coordinate) $\hat{\mathbf{v}}_c(t)$ for each frame. Similarly, we can find the hand contact vertex $\hat{\mathbf{h}}_c(t)$. During our experiments we find that it is sufficient to use the vertex located at the center of the palm. To make the curve matching pose-invariant and only focus on shape, we allow rigid transformation of the contact curve during matching as this does not change the shape of the curve. The pose-invariant human contact error is defined as $E_{\text{contact}}(\Theta_o^t, R_c; \boldsymbol{x}_h^t, H_{\text{est}}^t) =$

$$\sum_{t=t_{\text{start}}}^{t_{\text{end}}} (R_c\hat{\mathbf{v}}_c(t) + t_c - \hat{\mathbf{h}}_c(t))^2. \quad (6)$$

Here $R_c$ is the curve rotation and $t_c$ is the curve translation. In practice, we can directly compute $t_c$ as the offset between the first frames. Note that the rigid transformation is not frame-dependent.

### 7.4. Regularization

We use several regularization to avoid unrealistic results. The final regularization term $E_r$ includes 1) mask center error for pulling object mask closer to the annotated mask, 2) part motion error for penalizing out of limit motions, 3) depth error for making human and object depth closer during interaction, 4) global orientation error that penalizes large object roll or large negative pitch values (camera located beneath object), and finally 5) smoothing error that avoids sudden part motions. All these terms are formulated using the L2 loss.

### 8. Optimization

We perform gradient-based optimization with the objective goal: arg min $E(\Theta_o^t, \boldsymbol{x}_h^t, R_c)$ as defined in Eq. 2. Without assuming more information about the contact state, we optimize all possible combinations consisting of left/right hand, all object models from the corresponding category, and all possible object contact vertex locations. After optimization, the combination with the lowest $E(\Theta_o^t, \boldsymbol{x}_h^t, R_c)$ is chosen as our prediction. Adam [13] is used in all optimizations. The initial learning rate is set to $0.05$ and reduced to $0.005$ in the last 50 iterations. $\beta_1, \beta_2$ in Adam are set to $0.9, 0.999$ respectively. We optimize for a total of 200 iterations.

# 9. Experimental Results

We first discuss evaluation metrics for benchmarking our dataset in Section 9.1. Our optimization results are available in Section 9.2. We conduct ablation studies in Section 9.3. Finally we provide qualitative results in Section 9.4.

## 9.1. Evaluation Metrics

In our experiments, we use five different metrics for evaluating the pose, part motion and shape of the reconstructed object. We measure the difference in object orientation as the relative angle (in degrees) between estimated rotation matrix and the ground-truth rotation matrix (computed from Euler angles) in the SO(3) space. Object location error is measured as the 3D distance between estimated and ground-truth object base center (in cm). We calculate part motion error using absolute difference in degrees for revolute motion and absolute difference in length (cm) for prismatic motion. We also measure the averaged object dimension difference (in cm). Results in Section 9.2 are optimized with the ground truth CAD model being provided as part of the input. In the ablation studies, we also conduct studies where different models are considered and find the best matching one. In this case, we also report the CAD model correctness accuracy.

|  | orientation (degree) | location (cm) | part motion (degree/cm) | dimension (cm) |
|---|---|---|---|---|
| dishwasher | 10.134 | 20.941 | 12.158 | 12.287 |
| washingmachine | 19.791 | 28.673 | 19.634 | 12.243 |
| refrigerator | 11.301 | 31.390 | 13.343 | 11.013 |
| microwave | 22.504 | 18.901 | 14.560 | 10.991 |
| laptop | 29.379 | 29.535 | 23.530 | 8.469 |
| trashcan | 30.738 | 83.530 | 14.890 | 20.969 |
| oven | 7.059 | 14.644 | 14.273 | 12.788 |
| storage (revolute) | 12.036 | 16.542 | 17.486 | 6.860 |
| storage (prismatic) | 19.252 | 23.785 | 3.402 | 14.210 |
| average | 18.021 | 29.782 | 16.234 | 13.729 |

Table 2: Object state estimation error for full human-object interaction objective. Category-level error is reported using averaged result over all frames within one category. Last row provides average over all categories. Averaged part motion excludes storage (prismatic) due to its values been in cm rather than degrees.

## 9.2. Quantitative Evaluation

In this section, we provide quantitative results for our optimization-based method. We simplify the problem by using the ground truth CAD model. Note that we still need to estimate the object dimension. We use the ground-truth object masks (rendered from ground-truth object pose and motion provided in the dataset). Results for more challenging cases are examined in the ablation studies.

Table 2 summarizes the results on all eight categories. We separately report storage-furniture results based on

whether the moving part exhibits revolute or prismatic motion. We make several observations from the results. First, we see that part motion error is low in cases where small motion changes lead to large differences in the object mask. The part motion error for trashcans is the largest because the trashcan lid is very small and can not be easily reflected in our differentiable object mask term. Trashcans also have large location error because some of the captured trashcan data can not be correctly approximated well by the CAD models we use. Human pose estimation from Joo et al. [10] can also be noisy due to occlusions (e.g. laptop data where the laptop occludes most of the human body). The dimension error for large objects like refrigerators can also be very small. This indicates that it is not the case that large objects will always have larger size errors.

## 9.3. Ablation Studies

**Analysis of HOI term** To understand the effect of human-object relations, we compare results when optimized without human-object interaction error. The regularization term $E_\mathrm{r}$ is still used except for the human-object depth error. Results are averaged over all four best-performing categories (dishwasher, oven, storage revolute, and washingmachine). Figure 5 summarizes this comparison. We observe that optimized results deteriorate significantly without human-object interaction terms. The most decrease comes from global object orientation. This is expected since the 2D object mask does not contain much useful information to constrain the object pose.

**Analysis of CAD model** Previously, we assumed that the ground truth CAD model is given. In practice, we do not have access to this information. Thus it is necessary to optimize the selection of CAD model together with the rest of the parameters. To this end, we optimize over all CAD models in the corresponding category and keep the setting with the lowest objective cost. We test this on the dishwasher and storage (revolute joint) categories. The CAD model correctness accuracy is 90.91% and 71.43% for the two categories, respectively. When the retrieved model is incorrect, other errors will be large. This demonstrates the need for classification or retrieval approaches that can find the closest-matching CAD models.

## 9.4. Qualitative Evaluation

Figure 7 illustrates several example outputs after our optimization. Example videos of the results are available here. We observe that even though the human mesh is fixed in that we do not optimize its shape or joint angles, the human and object layout is visually plausible. The object dimensions ratio also look plausible and the object part moves accordingly as the human interacts with it.
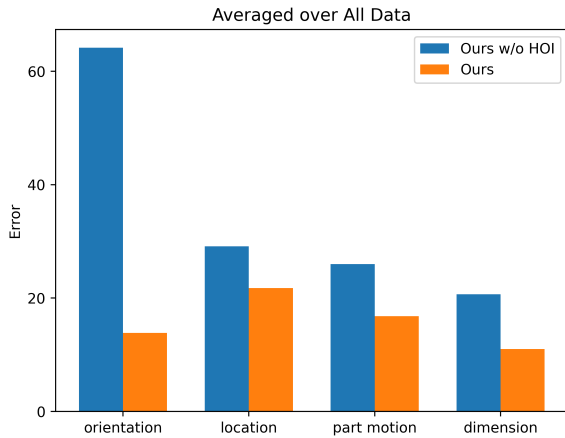
Figure 5: Comparison of full human-object interaction objective against object-only objective in terms of overall articulating part estimation error metrics for dishwasher, oven, storage (revolute), and washingmachine categories. Error is in degrees for orientation and part motion, and in cm for location and dimension. Lower is better. We see that removing the HOI terms significantly decreases performance.

We also visualize the optimization results without the HOI term. Results are shown in Figure 6. Without human-object relations, we find that the object orientation is incorrect in many cases. To still match the object mask at each frame, the estimated part motion is rather inaccurate. This either leads to large sudden motions or no motion. The same goes for the object dimension where the height, width, and length ratio is incorrect (e.g. oven).
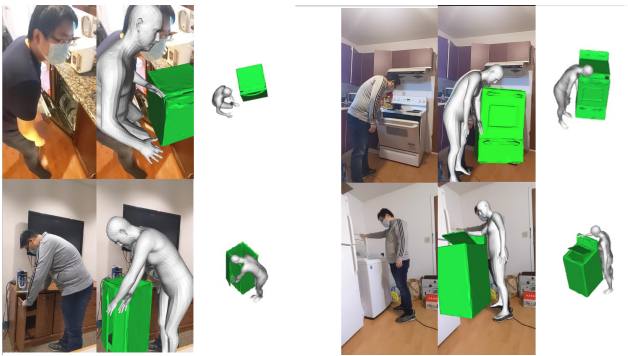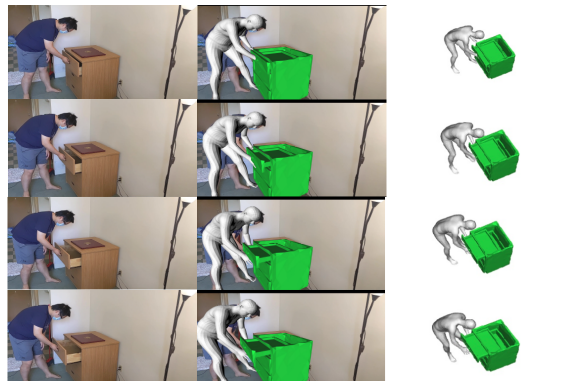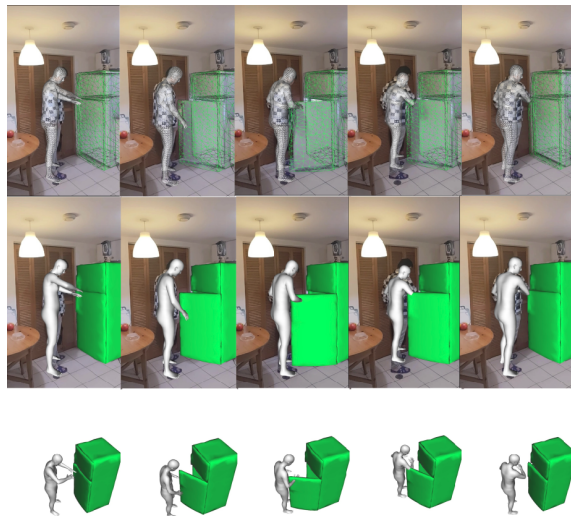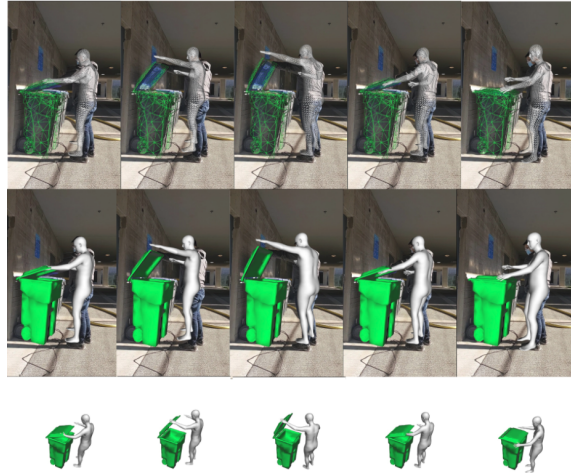


Figure 6: Results without the HOI term.



Figure 7: Results using the full human-object interaction objective optimization. The human and object layout is visually correct.

## 10. Conclusion

We introduce the D3D-HOI Dataset with frame-level annotation for object shape, pose, and part motion. Our dataset provides a real-world 3D benchmark for object pose and motion estimation during human object interaction. It can also be extended for training generative model for human-object interaction snippets. We propose

an optimization-based methods and show that modelling human-object relations can improve the estimated 3D object parameters. We hope this work will motivate further research in 3D human-object interactions.

## 11. Acknowledgement

## References

[1] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *European Conference on Computer Vision*, pages 561–578. Springer, 2016. 1, 2

[2] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000. 4

[3] Yixin Chen, Siyuan Huang, Tao Yuan, Siyuan Qi, Yixin Zhu, and Song-Chun Zhu. Holistic++ scene understanding: Single-view 3d holistic scene parsing and human pose estimation with human-object interaction and physical commonsense. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 2

[4] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9785–9795, 2019. 1, 2

[5] Mohamed Hassan, Vasileios Choutas, Dimitrios Tzionas, and Michael J Black. Resolving 3d human pose ambiguities with 3d scene constraints. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2282–2292, 2019. 1, 2

[6] Mohamed Hassan, Partha Ghosh, Joachim Tesch, Dimitrios Tzionas, and Michael J Black. Populating 3d scenes by learning human-scene interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14708–14718, 2021. 2

[7] Jiahui Huang, He Wang, Tolga Birdal, Minhyuk Sung, Federica Arrigoni, Shi-Min Hu, and Leonidas J Guibas. Multibodysync: Multi-body segmentation and motion estimation via 3d scan synchronization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7108–7118, 2021. 1, 2

[8] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325–1339, July 2014. 2

[9] Ajinkya Jain, Rudolf Lioutikov, and Scott Niekum. Screwnet: Category-independent articulation model estimation from depth images using screw theory. *arXiv preprint arXiv:2008.10518*, 2020. 2

[10] Hanbyul Joo, Natalia Neverova, and Andrea Vedaldi. Exemplar fine-tuning for 3d human pose fitting towards in-the-wild 3d human pose estimation. *arXiv*, 2020. 1, 2, 3, 5, 7

[11] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 1, 2

[12] Korrawe Karunratanakul, Jinlong Yang, Yan Zhang, Michael J Black, Krikamol Muandet, and Siyu Tang. Grasping field: Learning implicit representations for human grasps. In *2020 International Conference on 3D Vision (3DV)*, pages 333–344. IEEE, 2020. 2

[13] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6

[14] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. Pointrend: Image segmentation as rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9799–9808, 2020. 5

[15] Muhammed Kocabas, Nikos Athanasiou, and Michael J. Black. Vibe: Video inference for human body pose and shape estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1, 2

[16] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2252–2261, 2019. 1, 2

[17] Nilesh Kulkarni, Ishan Misra, Shubham Tulsiani, and Abhinav Gupta. 3d-relnet: Joint object and relational network for 3d prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2212–2221, 2019. 2

[18] Xiaolong Li, He Wang, Li Yi, Leonidas J Guibas, A Lynn Abbott, and Shuran Song. Category-level articulated object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3706–3715, 2020. 1, 2

[19] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7708–7717, 2019. 5

[20] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015. 2

[21] Zhengyi Luo, S. Alireza Golestaneh, and Kris M. Kitani. 3d human motion estimation via motion compression and refinement. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*, November 2020. 2

[22] Kevis-Kokitsi Maninis, Stefan Popov, Matthias Nießner, and Vittorio Ferrari. Vid2cad: Cad model alignment using multi-view constraints from videos. *arXiv preprint arXiv:2012.04641*, 2020. 1

[23] T. V. Marcard, Roberto Henschel, Michael J. Black, B. Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *ECCV*, 2018. 2

[24] Yinyu Nie, Xiaoguang Han, Shihui Guo, Yujian Zheng, Jian Chang, and Jian Jun Zhang. Total3dunderstanding: Joint layout, object pose and mesh reconstruction for indoor scenes from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages

55–64, 2020. 2

[25] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, P. Gehler, and B. Schiele. Neural body fitting: Unifying deep learning and model based human pose and shape estimation. *2018 International Conference on 3D Vision (3DV)*, pages 484–494, 2018. 2

[26] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 459–468, 2018. 2

[27] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10975–10985, 2019. 1, 2

[28] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 5

[29] Manolis Savva, Angel X Chang, Pat Hanrahan, Matthew Fisher, and Matthias Nießner. Pigraphs: learning interaction snapshots from observations. *ACM Transactions on Graphics (TOG)*, 35(4):1–12, 2016. 1, 2

[30] Bokui Shen, Fei Xia, Chengshu Li, Roberto Martín-Martín, Linxi Fan, Guanzhi Wang, Shyamal Buch, Claudia D'Arpino, Sanjana Srivastava, Lyne P Tchapmi, et al. igibson, a simulation environment for interactive tasks in large realisticscenes. *arXiv preprint arXiv:2012.02924*, 2020. 2

[31] Shubham Tulsiani, Saurabh Gupta, David F Fouhey, Alexei A Efros, and Jitendra Malik. Factoring shape, pose, and layout from the 2d image of a 3d scene. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 302–310, 2018. 1, 2

[32] Zhenzhen Weng and Serena Yeung. Holistic 3d human and scene mesh estimation from single view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 334–343, 2021. 1

[33] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, et al. Sapien: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11097–11107, 2020. 1, 2, 3, 4

[34] Zihao Yan, Ruizhen Hu, Xingguang Yan, Luanmin Chen, Oliver Van Kaick, Hao Zhang, and Hui Huang. Rpm-net: recurrent prediction of motion and parts from point cloud. *arXiv preprint arXiv:2006.14865*, 2020. 2

[35] Jason Y Zhang, Sam Pepose, Hanbyul Joo, Deva Ramanan, Jitendra Malik, and Angjoo Kanazawa. Perceiving 3d human-object spatial arrangements from a single image in the wild. In *European Conference on Computer Vision*, pages 34–51. Springer, 2020. 1, 2

[36] Siwei Zhang, Yan Zhang, Qianli Ma, Michael J. Black, and Siyu Tang. PLACE: Proximity learning of articulation and contact in 3D environments. In *International Conference on 3D Vision (3DV)*, November 2020. 2

[37] Yan Zhang, Mohamed Hassan, Heiko Neumann, Michael J Black, and Siyu Tang. Generating 3d people in scenes without people. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6194–6204, 2020. 2