

GRAB: A Dataset of Whole-Body Human Grasping of Objects

Omid Taheri, Nima Ghorbani, Michael J. Black, and Dimitrios Tzionas

Max Planck Institute for Intelligent Systems, Tübingen, Germany
{otaheri,nghorbani,black,dtzionas}@tuebingen.mpg.de

Abstract. Training computers to understand, model, and synthesize human grasping requires a rich dataset containing complex 3D object shapes, detailed contact information, hand pose and shape, and the 3D body motion over time. While “grasping” is commonly thought of as a single hand stably lifting an object, we capture the motion of the entire body and adopt the generalized notion of “whole-body grasps”. Thus, we collect a new dataset, called *GRAB* (GRasping Actions with Bodies), of whole-body grasps, containing full 3D shape and pose sequences of 10 subjects interacting with 51 everyday objects of varying shape and size. Given MoCap markers, we fit the full 3D body shape and pose, including the articulated face and hands, as well as the 3D object pose. This gives detailed 3D meshes over time, from which we compute contact between the body and object. This is a unique dataset, that goes well beyond existing ones for modeling and understanding how humans grasp and manipulate objects, how their full body is involved, and how interaction varies with the task. We illustrate the practical value of GRAB with an example application; we train GrabNet, a conditional generative network, to predict 3D hand grasps for unseen 3D object shapes. The dataset and code are available for research purposes at <https://grab.is.tue.mpg.de>.

1 Introduction

A key goal of computer vision is to estimate human-object interactions from video to help understand human behavior. Doing so requires a strong model of such interactions and learning this model requires data. However, capturing such data is not simple. Grasping involves both gross and subtle motions, as humans involve their whole body and dexterous finger motion to manipulate objects. Therefore, objects contact multiple body parts and not just the hands. This is difficult to capture with images because the regions of contact are occluded. Pressure sensors or other physical instrumentation, however, are also not a full solution as they can impair natural human-object interaction and do not capture full-body motion. Consequently, there are no existing datasets of complex human-object interaction that contain full-body motion, 3D body shape, and detailed body-object contact. To fill this gap, we capture a novel dataset of full-body 3D humans dynamically interacting with 3D objects as illustrated in Fig. 1. By accurately tracking 3D body and object shape, we reason about

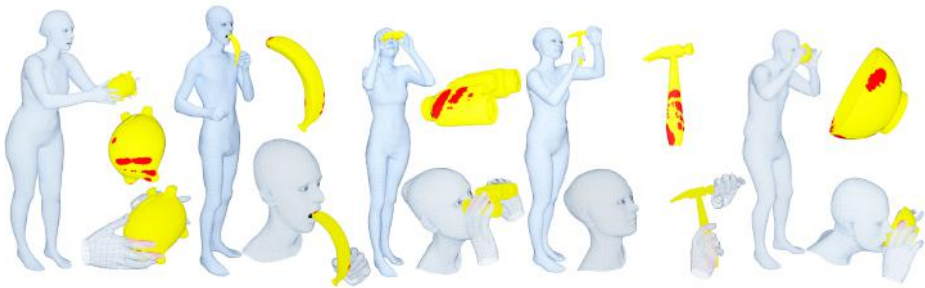


Fig. 1: Example “whole-body grasps” from the GRAB dataset. A “grasp” is usually thought of as a single hand interacting with an object. Using objects, however, may involve more than just a single hand. From left to right: (i) passing a piggy bank, (ii) eating a banana, (iii) looking through binoculars, (iv) using a hammer, (v) drinking from a bowl. Contact between the object and the body is shown in red on the object; here contact areas are spatially extended to aid visualization. See the video on our website for a wide range of sequences with various objects and intents.

contact resulting in a dataset with detail and richness beyond existing grasping datasets.

Most previous work focuses on prehensile “grasps” [47]; i.e. a single human hand stably lifting or using an object. The hands, however, are only part of the story. For example, as infants, our earliest grasps involve bringing objects to the mouth [62]. Consider the example of drinking from a bowl in Fig. 1 (right). To do so, we must pose our body so that we can reach it, we orient our head to see it, we move our arm and hand to stably lift it, and then we bring it to our mouth, making contact with the lips, and finally we tilt the head to drink. As this and other examples in the figure illustrate, human grasping and using of everyday objects involves the *whole body*. Such interactions are fundamentally *three-dimensional*, and contact occurs between objects and multiple body parts.

Dataset. Such whole-body grasping [28] has received much less attention [5,40] than single hand-object grasping [3,12,16,30,47]. To model such grasping we need a dataset of humans interacting with varied objects, capturing the full 3D surface of both the body and objects. To solve this problem we adapt recent motion capture techniques, to construct a new rich dataset called **GRAB** for “*G*Rasping *A*ctions with *B*odies.” Specifically, we adapt MoSh++ [41] in two ways. First, MoSh++ estimates the 3D shape and motion of the body and hands from MoCap markers; here we extend this to include facial motion. For increased accuracy we first capture a 3D scan of each subject and fit the SMPL-X body model [50] to it. Then MoSh++ is used to recover the pose of the body, hands and face. Note that the face is important because it is involved in many interactions; see in Fig. 1 (second from left) how the mouth opens to eat a banana. Second, we also accurately capture the motion of 3D objects as they are manipulated by the subjects. To this end, we use small hemi-spherical markers on

the objects and show that these do not impact grasping behavior. As a result, we obtain detailed 3D meshes for both the object and the human (with a full body, articulated fingers and face) moving over time while in interaction, as shown in Fig. 1. Using these meshes we then infer the body-object contact (red regions in Fig. 1). Unlike Brahmabhatt et al. [6] this gives both the contact and the full body/hand pose over time. Interaction is dynamic, including in-hand manipulation and re-grasping. GRAB captures 10 different people (5 male and 5 female) interacting with 51 objects from [6]. Interaction takes place in 4 different contexts: lifting, handing over, passing from one hand to the other, and using, depending on the affordances and functionality of the object.

Applications. GRAB supports multiple uses of interest to the community. First, we show how GRAB can be used to gain insights into hand-object contact in everyday scenarios. Second, there is a significant interest in training models to grasp 3D objects [4,63]. Thus, we use GRAB to train a pair of neural networks (coarse prediction followed by refinement) to generate plausible grasps for unseen 3D objects. Given a randomly posed 3D object, we predict plausible hand parameters (6 DoF wrist pose and full finger articulation) appropriate for grasping the object. To encode arbitrary 3D object shapes, we employ the recent basis point set (BPS) representation [56], whose fixed size is appropriate for neural networks. Then, by conditioning on a new 3D object shape, we sample from the learned latent space, and generate hand grasps for this. We quantitatively and qualitatively evaluate the resulting grasps and show that they look natural.

In summary, this work makes the following contributions: (1) we introduce a unique dataset capturing real “whole-body grasps” of 3D objects, including full-body human motion, object motion, in-hand manipulation and re-grasps; (2) to capture this, we adapt MoSh++ to solve for the body, face and hands of SMPL-X to obtain detailed moving 3D meshes; (3) using these meshes and tracked 3D objects we compute plausible contact on the object and the human and provide an analysis of observed patterns; (4) we show the value of our dataset for machine learning, by training a novel conditional neural network to generate 3D hand grasps for unseen 3D objects. The dataset, models, and code are available for research purposes at <https://grab.is.tue.mpg.de>.

2 Related Work

Hand Grasps: Hands are crucial for grasping and manipulating objects. For this reason, many studies focus on understanding grasps and defining taxonomies [3,12,16,30,47,55], by exploring the object shape and purpose of grasps [12], the contact areas on the hand captured by sinking objects in ink [30], the pose and contact areas [3] captured with an integrated data-glove [13] and tactile-glove [55], or the number of fingers in contact with the object and thumb position [16]. A key element for these studies is capturing accurate hand poses, relative hand-object configurations and contact areas.

Whole-Body Grasps: Often people use more than a single hand to interact with objects. However, there are not many works in the literature on this topic

[5,28]. Borras et al. [5] use MoCap data [42] of people interacting with a scene with multi-contact, and present a body pose taxonomy for such whole-body grasps. Hsiao et al. [28] focus on imitation learning with a database of whole-body grasp demonstrations with a human teleoperating a simulated robot. Although these works go in the right direction, they use unrealistic humanoid models and simple objects [5,28] or synthetic ones [28]. Instead, we use the SMPL-X model [50] to capture “whole-body”, face and dexterous in-hand interactions.

Capturing Interactions with MoCap: MoCap is often used to capture, synthesize or evaluate humans interacting with scenes. Lee et al. [39] capture a 3D body skeleton interacting with a 3D scene and show how to synthesize new motions in new scenes. Wang et al. [77] capture a 3D body skeleton interacting with large geometric objects. Han et al. [23] present a method for automatic labeling of hand markers, to speed up hand tracking for VR. Le et al. [38] capture a hand interacting with a phone to study the “comfortable areas”, while Feit et al. [15] capture two hands interacting with a keyboard to study typing patterns. Other works [37,53] focus on graphics applications. Kry et al. [37] capture a hand interacting with a 3D shape primitive, instrumented with a force sensor. Pollard et al. [53] capture the motion of a hand to learn a controller for physically based grasping. Mandery et al. [42] sit between the above works, capturing humans interacting with both big and handheld objects, but without articulated faces and fingers. None of the previous work captures full 3D bodies, hands and faces together with 3D object manipulation and contact.

Capturing Contact: Capturing human-object contact is hard, because the human and object heavily occlude each other. One approach is instrumentation with touch and pressure sensors, but this might bias natural grasps. Pham et al. [51] predefine contact points on objects to place force transducers. Recent advances in tactile sensors allow accurate recognition of tactile patterns and handheld objects [69]. Some approaches [3] use a data glove [13] with an embedded tactile glove [55,71], but this combination is complicated and the two modalities can be hard to synchronize. A microscopic-domain tactile sensor [19] is introduced in [29], but is not easy to use on human hands. Mascaro et al. [43] attach a minimally invasive camera to detect changes in the coloration of fingernails. Brahmabhatt et al. [6] use a thermal camera to directly observe the “thermal print” of a hand on the grasped object. However, for this they only capture static grasps that last long enough for heat transfer. Consequently, even recent datasets that capture realistic hand-object [18,22] or body-scene [26,65] interaction avoid directly measuring contact.

3D Interaction Models: Learning a model of human-object interactions is useful for graphics and robotics to help avatars [14,17,67] or robots [24] interact with their surroundings, and for vision [11,21,45] to help reconstruct interactions from ambiguous data. However, this is a chicken-and-egg problem; to capture or synthesize data to learn a model, one needs such a model in the first place. For this reason, the community has long used hand-crafted approaches that exploit contact and physics, for body-scene [7,25,26,61,80], body-object [33,35,40], or hand-object [27,49,58,59,66,72,73,76,81,82] scenarios. These approaches compute

contact approximately. This approximation may be rough when humans are modeled as 3D skeletons [33,40] or shape primitives [7,35,66]. However, it gets relatively accurate when using 3D meshes that are generic [58,72], personalized [25,61,73,76], or based on 3D statistical models [26,27].

To collect training data, several works [58,59] use synthetic Poser [54] hand models, manually articulated to grasp 3D shape primitives. Contact points and forces are also annotated [58] through proximity and inter-penetration of 3D meshes. In contrast, Hasson et al. [27] use the robotics method GraspIt [44] to automatically generate 3D MANO [60] grasps for ShapeNet [8] objects and render synthetic images of the hand-object interaction. However, GraspIt optimizes for hand-crafted grasp metrics that do not necessarily reflect the distribution of human grasps (see Sup. Mat. Sec. C.2 of [27], and [20]). Alternatively, Garcia-Hernando et al. [18] use magnetic sensors to reconstruct a 3D hand skeleton and rigid object poses; they capture 6 subjects interacting with 4 objects. This data is used for 3D hand/object pose estimation [36,70] or motion generation [17], but suffers from noisy poses and severe inter-penetrations (see Sec. 5.2 of [27]).

For bodies, Kim et al. [33] use synthetic data to learn to detect contact points on a 3D object, and then fit an interacting 3D body skeleton to them. Savva et al. [65] use RGB-D to capture 3D body skeletons of 5 subjects interacting in 30 3D scenes, to learn to synthesize interactions [65], affordance detection [64], or to reconstruct interaction from videos [45]. Mandery et al. [42] use optical MoCap to capture 43 subjects interacting with 41 tracked objects, both large and small. This is similar to our effort but they do not capture fingers or 3D body shape, so cannot reason about contact. Corona et al. [10] use this dataset to learn context-aware body motion prediction. Starke et al. [67] use Xsens IMU sensors [79] to capture the main body of a subject interacting with large objects, and learn to synthesize avatar motion in virtual worlds. Hassan et al. [26] use RGB-D and 3D scene constraints to capture 20 humans as SMPL-X [50] meshes interacting with 12 static 3D scenes, but do not capture object manipulation. Zhang et al. [83] use this data to learn to generate 3D scene-aware humans.

We see that only parts of our problem have been studied. We draw inspiration from prior work, in particular [6,26,28,42]. We go beyond these by introducing a new dataset of real “whole-body” grasps, as described in the next section.

3 Dataset

To manipulate an object, the human needs to approach its 3D surface, and bring their skin to come in *physical contact* to apply forces. Realistically capturing such human-object interactions, especially with “whole-body grasps”, is a challenging problem. First, the object may occlude the body and vice-versa, resulting in *ambiguous* observations. Second, for physical interactions it is crucial to reconstruct an accurate and detailed 3D *surface* for both the human and the object. Additionally, the capture has to work across multiple scales (body, fingers and face) and for objects of varying complexity. We address these challenges with a unique combination of state-of-the-art solutions that we adapt to the problem.



Fig. 2: MoCap markers used to capture humans and objects. **Left:** We attach 99 reflective markers per subject; 49 for the body, 14 for the face and 36 for the fingers. We use spherical 4.5 mm radius markers for the body and hemi-spherical 1.5 mm radius ones for the hands and face. **Right:** Example 3D printed objects from [6]. We glue 1.5 mm radius hemi-spherical markers (the gray dots) on the objects. These markers are small enough to be unobtrusive. The 6 objects on the right are mostly used by one or more hands, while the 6 on the left involve “whole-body grasps”.

There is a fundamental trade-off with current technology; one has to choose between (a) accurate motion with instrumentation and without natural RGB images, or (b) less accurate motion but with RGB images. Here we take the former approach; for an extensive discussion we refer the reader to Sup. Mat.

3.1 Motion Capture (MoCap)

We use a Vicon system with 54 infrared “Vantage 16” [75] cameras that capture 16 MP at 120 fps. The large number of cameras minimizes occlusions and the high frame rate captures temporal details of contact. The high resolution allows to capture even small (1.5 mm radius) hemi-spherical markers. This minimizes their influence on finger and face motion and does not alter how people grasp objects. Details of the marker setup are shown in Fig. 2. Even with many cameras, motion capture of the body, face, and hands, together with objects, is uncommon because it is so challenging. MoCap markers become occluded, labels are swapped, and ghost makers appear. MoCap cleaning was done by four trained technicians using Vicon’s Shōgun-Post software.

Capturing Human MoCap: To capture human motion, we use the marker set of Fig. 2 (left). The body markers are attached on a tight body suit with velcro-based mounting at a distance of roughly $d_b = 9.5$ mm from the body surface. The hand and face markers are attached directly on the skin with special removable glue, therefore the distance to it is roughly $d_h = d_f \approx 0$ mm. Importantly, no hand glove is used and hand markers are placed only on the dorsal side, leaving the palmar side completely uninstrumented, for natural interactions.

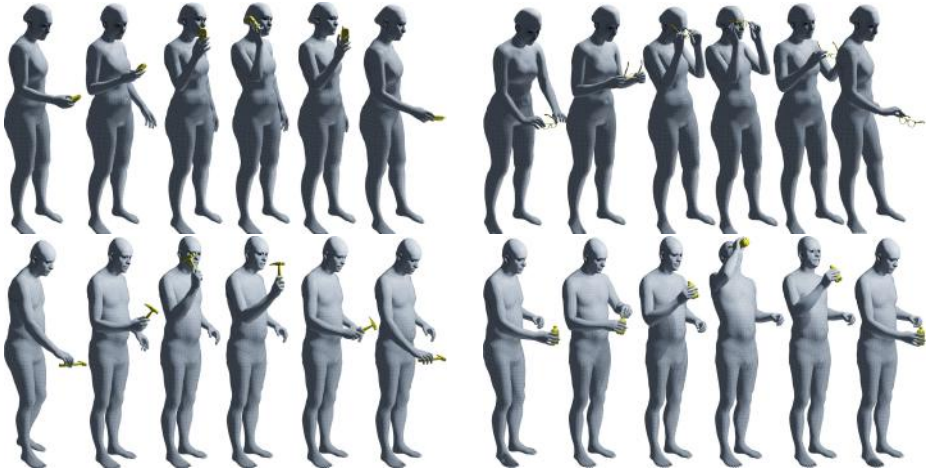


Fig. 3: We capture humans interacting with objects over time and reconstruct sequences of 3D meshes for both, as described in Sec. 3.1 and Sec. 3.2. Note the realistic and plausible placement of objects in the hands, and the “whole-body” involvement. The video on our website shows more examples.

Capturing Objects: To reconstruct interactions accurately, it is important to know the precise 3D object surface geometry. We therefore use the CAD object models of [6], and 3D print them with a Stratasys Fortus 360mc [68] printer; see Fig. 2 (right). Each object o is then represented by a known 3D mesh with vertices V_o . To capture object motion, we attach 1.5 mm radius hemi-spherical markers directly on the object surface with strong glue. We use at least 8 markers per object, empirically distributing them on the object so that at least 3 of them are always observed. The size and placement of the markers makes them unobtrusive. In Sup. Mat. we show empirical evidence that markers have minimal influence on grasping.

3.2 From MoCap Markers to 3D Surfaces

Human Model: We model the human with the SMPL-X [50] 3D body model. SMPL-X jointly models the body with an articulated face and fingers; this expressive body model is critical to capture physical interactions. More formally, SMPL-X is a differentiable function $M_b(\beta, \theta, \psi, \gamma)$ that is parameterized by body shape β , pose θ , facial expression ψ and translation γ . The output is a 3D mesh $M_b = (V_b, F_b)$ with $N_b = 10475$ vertices $V_b \in \mathbb{R}^{(N_b \times 3)}$ and triangles F_b . The shape parameters $\beta \in \mathbb{R}^{10}$ are coefficients in a learned low-dimensional linear shape space. This lets SMPL-X represent different subject identities with the same mesh topology. The 3D joints, $J(\beta)$, of a kinematic skeleton are regressed from the body shape defined by β . The skeleton has 55 joints in total; 22 for the body, 15 joints per hand for finger articulation, and 3 for the neck and eyes. Corrective blend shapes are added to the body shape, and then the posed body is

defined by linear blend skinning with this underlying skeleton. The overall pose parameters $\theta = (\theta_b, \theta_f, \theta_h)$ are comprised of $\theta_b \in \mathbb{R}^{66}$ and $\theta_f \in \mathbb{R}^9$ parameters in axis-angle representation for the main body and face joints correspondingly, with 3 degrees of freedom (DoF) per joint, and $\theta_h \in \mathbb{R}^{60}$ parameters in a lower-dimensional pose space for both hands, i.e. 30 DoF per hand following [27]. For more details, please see [50].

Model-Marker Correspondences: For the human body we define, a priori, the rough marker placement on the body as shown in Fig. 2 (left). Exact marker locations on individual subjects are then computed automatically using MoSh++ [41]. In contrast to the body, the objects have different shapes and mesh topologies. Markers are placed according to the object shape, affordances and expected occlusions during interaction; see Fig. 2 (right). Thus, we annotate object-specific vertex-marker correspondences, and do this once per object.

Human and Object Tracking: To ensure accurate human shape, we capture a 3D scan of each subject and fit SMPL-X to it following [60]. We fit these personalized SMPL-X models to our cleaned 3D marker observations using MoSh++ [41]. Specifically we optimize over pose, θ , expressions, ψ , and translation, γ , while keeping the known shape, β , fixed. The weights of MoSh++ for the finger and face data terms are tuned on a synthetic dataset, as described in Sup. Mat. An analysis of MoSh++ fitting accuracy is also provided in Sup. Mat.

Objects are simpler because they are rigid and we know their 3D shape. Given three or more detected markers, we solve for the rigid object pose $\theta_o \in \mathbb{R}^6$. Here we track the human and object separately and on a per-frame basis. Figure 3 shows that our approach captures realistic interactions and reconstructs detailed 3D meshes for both the human and the object, over time. The video on our website shows a wide range of reconstructed sequences.

3.3 Contact Annotation

Since contact cannot be directly observed, we estimate it using 3D proximity between the 3D human and object meshes. In theory, they come in contact when the distance between them is zero. In practice, however, we relax this and define contact when the distance, d , is below a tolerance, threshold $d \leq \epsilon_{\text{contact}}$. This helps address: (1) measurement and fitting errors, (2) limited mesh resolution, (3) the fact that human soft tissue deforms when grasping an object, while the SMPL-X model cannot model this.

Given these issues, accurately estimating contact is challenging. Consider the hand grasping a wine glass in Fig. 4 (right), where the color rings indicate intersections. Ideally, the glass should be in contact with the thumb, index and middle fingers. “Contact under-shooting” results in fingers hovering close to the object surface, but not on it, like the thumb. “Contact over-shooting”, results in fingers penetrating the object surface around the contact area, like the index (purple intersections) and middle finger (red intersections). The latter case is especially problematic for thin objects where a penetrating finger can pass through the object, intersecting it on two sides. In this example, we want to annotate contact only with the outer surface of the object and not the inner one.

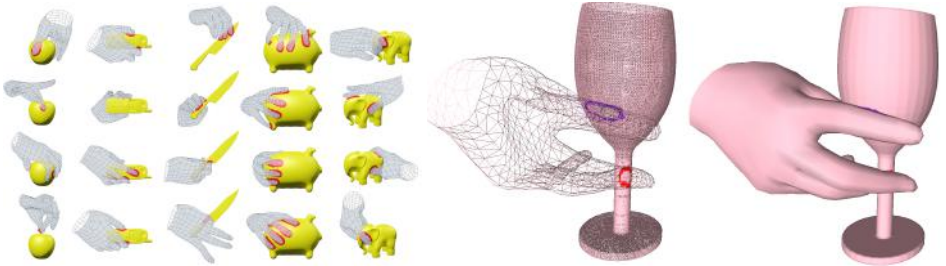


Fig. 4: Left: Accurate tracking lets us compute realistic contact areas (red) for each frame (Sec. 3.3). For illustration, we render only the hand of SMPL-X and spatially extend the red contact areas for visibility. **Right:** Detection of “intersection ring” triangles during contact annotation (Sec. 3.3).

We account for “contact over-shooting” cases with an efficient heuristic. We use a fast method [32,50] to detect intersections, cluster them in connected “intersection rings”, $\mathcal{R}_b \subsetneq V_b$ and $\mathcal{R}_o \subsetneq V_o$, and label them with the intersecting body part, seen as purple and red rings in Fig. 4 (right). The “intersection ring”, \mathcal{R}_b , segments the body mesh, M_b , to give the “penetrating sub-mesh”, $\mathcal{M}_b \subsetneq M_b$. We then identify two cases: (1) When a body part gives only one intersection, we annotate as contact points on the object, $V_o^C \subset V_o$, all vertices enclosed by the ring \mathcal{R}_o . We then annotate as contact points on the body, $V_b^C \subset V_b$, all vertices that lie close to V_o^C with a distance $d_{o \rightarrow b} \leq \epsilon_{contact}$. (2) In case of multiple intersections, i , we take into account only the ring \mathcal{R}_b^i corresponding to the largest intersection subset, \mathcal{M}_b^i .

For body parts that are not found in contact above, there is the possibility of “contact under-shooting”. To address this, we compute the distance from each object vertex, V_o , to each non-intersecting body vertex, V_b . We then annotate as contact vertices, V_o^C and V_b^C , the ones with $d_{o \rightarrow b} \leq \epsilon_{contact}$. We empirically find that $\epsilon_{contact} = 4.5$ mm works well for our purposes.

3.4 Dataset Protocol

Human-object interaction depends on various factors including the human body shape, object shape and affordances, object functionality, or interaction intent, to name a few. We therefore capture 10 people (5 men and 5 women), of various sizes and nationalities, interacting with the objects of [6]; see example objects in Fig. 2 (right). All subjects gave informed written consent to share their data for research purposes.

For each object we capture interactions with 4 different intents, namely “use” and “pass” (to someone), borrowed from the protocol of [6], as well as “lift” and “off-hand pass” (from one hand to the other). Figure 3 shows some example 3D capture sequences for the “use” intent. For each sequence we: (i) we randomize initial object placement to increase motion variance, (ii) we instruct the subject to follow an intent, (iii) the subject starts from a T-pose and approaches the

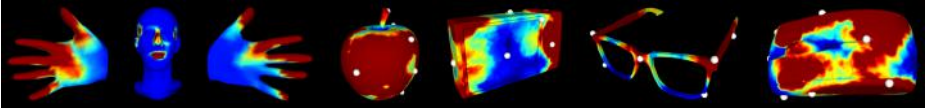


Fig. 5: Contact “heatmaps”. **Left:** For the body we focus on “use” sequences to show “whole-body grasps”. **Right:** For objects we include all intents. Object markers (light gray) are unobtrusive and can lie on “hot” (red) contact areas.

Table 1: Size of the GRAB dataset. GRAB is sufficiently large to enable training of data-driven models of grasping as shown in Sec. 4.

Intent	“Use”	“Pass”	“Lift”	“Off-hand”	Total
# Sequences	579	414	274	67	1334
# Frames	605.796	335.733	603.381	77.549	1.622.459

object, (iv) performs the instructed task, and (v) leaves the object and returns to a T-pose. The video on our website shows the richness of our protocol with a wide range of captured sequences.

3.5 Analysis

Dataset Analysis. The dataset contains 1334 sequences and over 1.6M frames of MoCap; Table 1 provides a detailed breakdown. Here we analyze those frames where we have detected contact between the human and the object. We assume that the object is static on a table and can move only due to grasping. Consequently, we consider contact frames to be those in which the object’s position deviates in the vertical direction by at least 5 mm from its initial position and in which at least 50 body vertices are in contact with the object. This results in 952,514 contact frames that we analyze below. The exact thresholds of these contact heuristics have little influence on our analysis, see Sup. Mat.

By uniquely capturing the whole body, and not just the hand, interesting interaction patterns arise. By focusing on “use” sequences that highlight the object functionality, we observe that 92% of contact frames involve the right hand, 39% the left hand, 31% both hands, and 8% involve the head. For the first category the per-finger contact likelihood, from thumb to pinky, is 100%, 96%, 92%, 79%, 39% and for the palm 24%. For more results see Sup. Mat.

To visualize the fine-grained contact information, we integrate over time the binary per-frame contact maps, and generate “heatmaps” encoding the contact likelihood across the whole body surface. Figure 5 (left) shows such “heatmaps” for “use” sequences. “Hot” areas (red) denote high likelihood of contact, while “cold” areas (blue) denote low likelihood. We see that both the hands and face are important for using everyday objects, highlighting the importance of capturing the whole interacting body. For the face, the “hot” areas are the lips, the nose, the temporal head area, and the ears. For hands, the fingers are more frequently in contact than the palm, with more contact on the right hand than

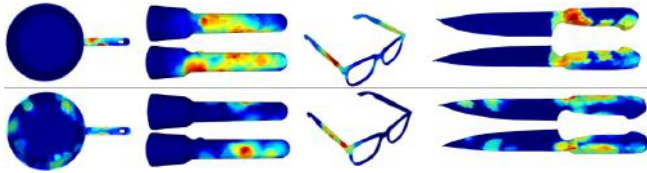


Fig. 6: Effect of interaction intent on contact during grasping. We show the “use” (top) and “pass” (bottom) intents for 4 different objects.

the left one. The palm seems more important for right-hand grasps than for left-hand ones, possibly because all our subjects are right-handed. Contact patterns are also influenced by the size of the object and the size of the hand; see Sup. Mat. for a visualization.

Figure 6 shows the effect of the intent. Contact for “use” sequences complies with the functionality of the object; e.g. people do not touch the knife blade or the hot area of the pan, but they do contact the on/off button of the flashlight. For “pass” sequences subjects tend to contact one side of the object irrespective of affordances, leaving the other one free to be grasped by the receiving person.

For natural interactions it is important to have a minimally intrusive setup. While our MoCap markers are small and unobtrusive, as seen in Figure 2 (right), we ask whether subjects may be biased in their grasps by these markers. Figure 5 (right) shows contact “heatmaps” for some objects across all intents. These clearly show that markers are often located in “hot” areas, suggesting that subjects do not avoid grasping these locations. Further analysis based on K-means clustering of grasps can be found in Sup. Mat.

4 GrabNet: Learning to Grab an Object

We show the value of GRAB with a challenging application; we train on it a model that generates plausible 3D MANO [60] grasps for an unseen 3D object. Our model, GrabNet, is comprised of two modules: coarse prediction and refinement. This is similar to [46], but with several key differences. We predict full human hand pose instead of a robotic gripper, train using captured human grasps, employ a different object representation which generalizes to new objects, and make different choices for the model and losses. Importantly, our refinement is done by a neural network, not an optimization process.

We first employ CoarseNet, a conditional variational autoencoder (cVAE) [34], that generates an initial grasp. For this it learns a grasping embedding space, Z , conditioned on the object shape, that is encoded using the Basis Point Set (BPS) [56] representation as a set of distances from the basis points to the nearest object points. In contrast to [46], Z captures not only the 6 DoF pre-grasp pose (gripper for [46]/wrist for MANO), but also the fully articulated human hand pose. CoarseNet’s grasps are reasonable, but realism can improve by refining contacts based on the distances, D , between the hand and object

meshes. We do this with a second network, called RefineNet, that employs explicit MANO contact points learned on our GRAB dataset to refine the initial pose. In contrast, [46] learn an evaluator of the 6 DoF gripper pose that they differentiate to optimize the coarse pose. The architecture of GrabNet is shown in Fig. 7, for more details see Sup. Mat.

Pre-processing. For training, we gather all frames with right-hand grasps that involve some minimal contact, for details see Sup. Mat. We then center each training sample, i.e. hand-object grasp, at the centroid of the object and compute the $BPS_o \in R^{4096}$ representation for the object, used for conditioning.

CoarseNet. We pass the object shape BPS_o along with initial MANO wrist rotation θ_{wrist} and translation γ to the encoder $Q(Z|\theta_{wrist}, \gamma, BPS_o)$ that produces a latent grasp code $Z \in R^{16}$. The decoder $P(\hat{\theta}, \hat{\gamma}|Z, BPS_o)$ maps Z and BPS_o to MANO parameters with full finger articulation $\hat{\theta}$, to generate a 3D grasping hand. For the training loss, we use standard cVAE loss terms (KL divergence, weight regularizer), a data term on MANO mesh edges (L1), as well as a penetration and a contact loss. For the latter, we learn candidate contact point weights from GRAB, in contrast to handcrafted ones [26] or weights learned from artificial data [27]. At inference time, given an unseen object shape, BPS_o , we sample the latent space, Z , and decode our sample to generate a MANO grasp.

RefineNet. The grasps estimated by CoarseNet are plausible, but can be refined for improved contacts. For this, RefineNet takes as input the initial grasp ($\hat{\theta}$, $\hat{\gamma}$) and the distances D from MANO vertices to the object mesh. The distances are weighted according to the vertex contact likelihood learned from GRAB. Then, RefineNet estimates refined MANO parameters ($\hat{\theta}$, $\hat{\gamma}$) in 3 iterative steps as in [31], to give the final grasp. To train RefineNet, we generate a synthetic dataset; we sample CoarseNet grasps as ground truth and we perturb their hand pose parameters to simulate noisy input estimates. We use the same training losses as for CoarseNet.

GrabNet. Given an unseen 3D object, we first get an initial grasp estimate with CoarseNet, and pass this to RefineNet to get the final grasp estimate. For simplicity, the two networks are trained separately, but we expect end-to-end refinement to be beneficial, as in [27]. Figure 8 (right) shows generated examples; our generations look realistic, as explained later in the evaluation section. For more qualitative results, see the video on our website and images in Sup. Mat.

Contact. As a free by-product of our 3D grasp predictions, we can compute contact between the 3D hand and object meshes, following Sec. 3.3. Contacts for GrabNet are shown with red in Figure 8 (right). Other methods for contact prediction, like [6], are pure bottom-up approaches that label a vertex as in contact or not, without explicit reasoning about the hand structure. In contrast, we follow a top-down approach; we first generate a 3D grasping hand, and then compute contact with explicit anthropomorphic reasoning.

Evaluation - CoarseNet/RefineNet. We first quantitatively evaluate the two main components, by computing the reconstruction vertex-to-vertex error. For CoarseNet the errors are 12.1 mm, 14.1 mm and 18.4 mm for the training, validation and test set respectively. For RefineNet the errors are 3.7 mm, 4.1 mm

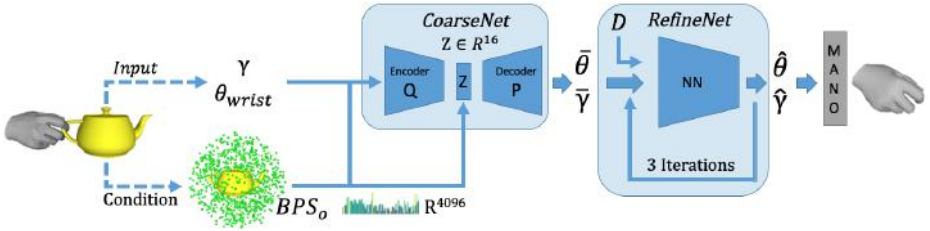


Fig. 7: GrabNet architecture. GrabNet generates MANO [60] grasps for unseen object shapes, encoded with a BPS [56] representation. It is comprised of two main modules. First, with CoarseNet we predict an initial plausible grasp. Second, we refine this with RefineNet to produce better contacts with the object.

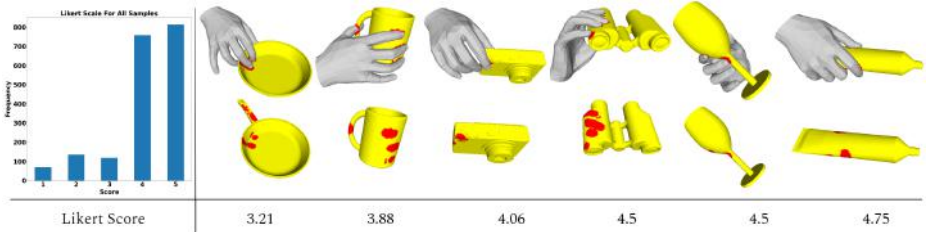


Fig. 8: Grasps generated by GrabNet for unseen objects; grasps look natural. As free by-product of 3D mesh generation, we get the red contact areas. For each grasp we show the average Likert score from all annotators. On the left we show the average Likert score for all generated grasps (best viewed on screen).

and 4.4 mm. The results show that the components, that are trained separately, work reasonably well before plugging them together.

Evaluation - GrabNet. To evaluate GrabNet generated grasps, we perform a user study through AMT [1]. We take 6 test objects from the dataset and, for each object, we generate 20 grasps, mix them with 20 ground-truth grasps, and show them with a rotating 3D viewpoint to subjects. Then we ask participants how they agree with the statement “Humans can grasp this object as the video shows” on a 5-level Likert scale (5 is “strongly agree” and 1 is “strongly disagree”). To filter out the noisy subjects, namely the ones who do not understand the task or give random answers, we use catch trials that show implausible grasps. We remove subjects who rate these catch trials as realistic; see Sup. Mat. for details. Table 2 (left) shows the user scores for both ground-truth and generated grasps.

Evaluation - Contact. Figure 9 shows examples of contact areas (red) generated by [6] (left) and our approach (right). The method of [6] gives only 10 predictions per object, some with zero contact. Also, a hand is supposed to touch the whole red area; this is often not anthropomorphically plausible. Our contact is a by-product of MANO-based inference and is, by construction, anthropomorphically valid. Also, one can draw infinite samples from our learned

Table 2: GrabNet evaluation for 6 test objects. The “AMT” column shows user study results; grasp quality is rated from 1 (worst) to 5 (best). The “vertices” and “contact” columns evaluate grasps against the closest ground-truth one.

Test Object	AMT				Vertices	
	Generation		Ground Truth		mean(cm)	Contact
	mean	std	mean	std	N=100	%
binoculars	4.09	0.93	4.27	0.80	2.56	4.00
camera	4.40	0.79	4.34	0.76	2.90	3.75
frying pan	3.19	1.30	4.49	0.67	3.58	4.16
mug	4.13	1.00	4.36	0.78	1.96	3.25
toothpaste	4.56	0.67	4.42	0.77	1.78	5.39
wineglass	4.32	0.88	4.43	0.79	1.92	4.56
Average	4.12	1.04	4.38	0.77	2.45	4.18

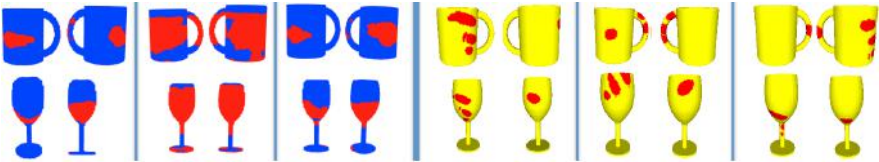


Fig. 9: Comparison of GrabNet (right) to ContactDB [6] (left). For each estimation we render two views, following the presentation style of [6].

grasping latent space. For further evaluation, we follow a protocol similar to [6] for our data. For every unseen test object we generate 20 grasps, and for each one we find both the closest ground-truth contact map and the closest ground-truth hand vertices, for comparison. Table 2 (right) reports the average error over all 20 predictions, in % for the former and cm for the latter case.

5 Discussion

We provide a new dataset to the community that goes beyond previous motion capture or grasping datasets. We believe that GRAB will be useful for a wide range of problems. Here we show that it provides enough data and variability to train a novel network to predict human grasping of objects, as we demonstrate with GrabNet. But there is much more that can be done. Importantly, GRAB includes the whole-body motion, enabling a much richer modeling than GrabNet.

Limitations: By focusing on accurate MoCap, we do not have synced image data. However, GRAB can support image-based inference [9,11,31] by enabling rendering of synthetic human-object interaction [27,57,74] or learning priors to regularize ill-posed inference of human-object interaction from 2D images [45].

Future Work: GRAB can support learning human-object interaction models [45,65], robotic grasping from imitation [17,78], mapping MoCap markers to meshes [23], rendering synthetic images [27,57,74], inferring object shape/pose from interaction [2,48], or analysis of temporal patterns [52].

Acknowledgements: We thank S. Polikovsky, M. Höschle (MH) and M. Landry (ML) for the MoCap facility. We thank F. Mattioni, D. Hieber, and A. Valis for MoCap cleaning. We thank ML and T. Alexiadis for trial coordination, MH and F. Grimmer for 3D printing, V. Callaghan for voice recordings and J. Tesch for renderings.

Disclosure: In the last five years, MJB has received research gift funds from Intel, Nvidia, Facebook, and Amazon. He is a co-founder and investor in Meshcapade GmbH, which commercializes 3D body shape technology. While MJB is a part-time employee of Amazon, his research was performed solely at, and funded solely by, MPI.

References

1. Amazon Mechanical Turk. <https://www.mturk.com>
2. Behbahani, F.M., Singla-Buxarraais, G., Faisal, A.A.: Haptic SLAM: An ideal observer model for bayesian inference of object shape and hand pose from contact dynamics. In: Haptics: Perception, Devices, Control, and Applications (2016)
3. Bernardin, K., Ogawara, K., Ikeuchi, K., Dillmann, R.: A sensor fusion approach for recognizing continuous human grasping sequences using hidden markov models. *IEEE Transactions on Robotics (T-RO)* **21**(1), 47–57 (2005)
4. Bohg, J., Morales, A., Asfour, T., Kragic, D.: Data-driven grasp synthesisa survey. *Transaction on Robotics (TRO)* **30**(2), 289309 (Apr 2014)
5. Borras, J., Asfour, T.: A whole-body pose taxonomy for loco-manipulation tasks. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (2015)
6. Brahmmbhatt, S., Ham, C., Kemp, C.C., Hays, J.: ContactDB: Analyzing and predicting grasp contact via thermal imaging. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
7. Brubaker, M.A., Fleet, D.J., Hertzmann, A.: Physics-based person tracking using the anthropomorphic walker. *International Journal of Computer Vision (IJCV)* **87**(1), 140 (2009)
8. Chang, A.X., tomas A. Funkhouser, eonidas J. Guibas, Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., Xiao, J., Yi, L., Yu, F.: ShapeNet: An information-rich 3D model repository. *arXiv:1512.03012* (2015)
9. Choutas, V., Pavlakos, G., Bolkart, T., Tzionas, D., Black, M.J.: Monocular expressive body regression through body-driven attention. In: *Proceedings of the European Conference on Computer Vision (ECCV)* (2020)
10. Corona, E., Pumarola, A., Alenyà, G., Moreno-Noguer, F.: Context-aware human motion prediction. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020)
11. Corona, E., Pumarola, A., Alenya, G., Moreno-Noguer, F., Rogez, G.: GanHand: Predicting human grasp affordances in multi-object scenes. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020)
12. Cutkosky, M.R.: On grasp choice, grasp models, and the design of hands for manufacturing tasks. *IEEE Transactions on Robotics and Automation* **5**(3), 269–279 (1989)
13. Cyberglove III data glove. <http://www.cyberglovesystems.com/cyberglove-iii>
14. ElKoura, G., Singh, K.: Handrix: animating the human hand. In: *Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation* (2003)

15. Feit, A.M., Weir, D., Oulasvirta, A.: How we type: Movement strategies and performance in everyday typing. In: *Proceedings of the CHI Conference on Human Factors in Computing Systems* (2016)
16. Feix, T., Romero, J., Schmiedmayer, H.B., Dollar, A.M., Kragic, D.: The GRASP taxonomy of human grasp types. *IEEE Transactions on Human-Machine Systems* **46**(1), 66–77 (2016)
17. Garcia-Hernando, G., Johns, E., Kim, T.K.: Physics-based dexterous manipulations with estimated hand poses and residual reinforcement learning. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (2020)
18. Garcia-Hernando, G., Yuan, S., Baek, S., Kim, T.K.: First-person hand action benchmark with rgb-d videos and 3d hand pose annotations. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2018)
19. GelSight tactile sensor. <http://www.gelsight.com>
20. Goldfeder, C., Ciocarlie, M.T., Dang, H., Allen, P.K.: The Columbia grasp database. In: *IEEE International Conference on Robotics and Automation (ICRA)* (2009)
21. Hamer, H., Gall, J., Weise, T., Van Gool, L.: An object-dependent hand pose prior from sparse training data. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2010)
22. Hampali, S., Oberweger, M., Rad, M., Lepetit, V.: HO-3D: A multi-user, multi-object dataset for joint 3D hand-object pose estimation. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
23. Han, S., Liu, B., Wang, R., Ye, Y., Twigg, C.D., Kin, K.: Online optical marker-based hand tracking with deep labels. *ACM Transactions on Graphics (TOG)* **37**(4), 166:1166:10 (2018)
24. Handa, A., Wyk, K.V., Yang, W., Liang, J., Chao, Y.W., Wan, Q., Birchfield, S., Ratliff, N., Fox, D.: DexPilot: Vision based teleoperation of dexterous robotic hand-arm system. In: *IEEE International Conference on Robotics and Automation (ICRA)* (2019)
25. Hasler, N., Rosenhahn, B., Thormahlen, T., Wand, M., Gall, J., Seidel, H.: Markerless motion capture with unsynchronized moving cameras. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2009)
26. Hassan, M., Choutas, V., Tzionas, D., Black, M.J.: Resolving 3D human pose ambiguities with 3D scene constraints. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2019)
27. Hasson, Y., Varol, G., Tzionas, D., Kalevatykh, I., Black, M.J., Laptev, I., Schmid, C.: Learning joint reconstruction of hands and manipulated objects. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
28. Hsiao, K., Lozano-Perez, T.: Imitation learning of whole-body grasps. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems* (2006)
29. Johnson, M.K., Cole, F., Raj, A., Adelson, E.H.: Microgeometry capture using an elastomeric sensor. *ACM Transactions on Graphics (TOG)* **30**(4), 46:1–46:8 (2011)
30. Kamakura, N., Matsuo, M., Ishii, H., Mitsuboshi, F., Miura, Y.: Patterns of static prehension in normal hands. *American Journal of Occupational Therapy* **34**(7), 437–445 (1980)
31. Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2018)
32. Karras, T.: Maximizing parallelism in the construction of bvhs, octrees, and k-d trees. In: *Proceedings of the ACM SIGGRAPH/Eurographics Conference on High-Performance Graphics* (2012)

33. Kim, V.G., Chaudhuri, S., Guibas, L., Funkhouser, T.: Shape2pose: Human-centric shape analysis. *ACM Transactions on Graphics (TOG)* **33**(4), 120:1–120:12 (2014)
34. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: *International Conference on Learning Representations (ICLR)* (2014)
35. Kjellstrom, H., Kragic, D., Black, M.J.: Tracking people interacting with objects. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2010)
36. Kovic, M., Kragic, D., Bohg, J.: Learning task-oriented grasping from human activity datasets. *IEEE Robotics and Automation Letters (RA-L)* **5**(2), 3352–3359 (2020)
37. Kry, P.G., Pai, D.K.: Interaction capture and synthesis. *ACM Transactions on Graphics (TOG)* **25**(3), 872–880 (2006)
38. Le, H.V., Mayer, S., Bader, P., Henze, N.: Fingers’ range and comfortable area for one-handed smartphone interaction beyond the touchscreen. In: *Proceedings of the CHI Conference on Human Factors in Computing Systems* (2018)
39. Lee, K.H., Choi, M.G., Lee, J.: Motion patches: building blocks for virtual environments annotated with motion data. *ACM Transactions on Graphics (TOG)* **25**(3), 898–906 (2006)
40. Li, Z., Sedlar, J., Carpentier, J., Laptev, I., Mansard, N., Sivic, J.: Estimating 3D motion and forces of person-object interactions from monocular video. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
41. Mahmood, N., Ghorbani, N., F. Troje, N., Pons-Moll, G., Black, M.J.: AMASS: Archive of motion capture as surface shapes. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2019)
42. Mandery, C., Terlemez, Ö., Do, M., Vahrenkamp, N., Asfour, T.: The KIT whole-body human motion database. In: *International Conference on Advanced Robotics (ICAR)* (2015)
43. Mascaro, S.A., Asada, H.H.: Photoplethysmograph fingernail sensors for measuring finger forces without haptic obstruction. *IEEE Transactions on Robotics and Automation (TRA)* **17**(5), 698–708 (2001)
44. Miller, A.T., Allen, P.K.: Graspit! A versatile simulator for robotic grasping. *IEEE Robotics Automation Magazine (RAM)* **11**(4), 110–122 (2004)
45. Monszpart, A., Guerrero, P., Ceylan, D., Yumer, E., Mitra, N.J.: iMapper: Interaction-guided scene mapping from monocular videos. *ACM Transactions on Graphics (TOG)* **38**(4), 92:1–92:15 (2019)
46. Mousavian, A., Eppner, C., Fox, D.: 6-dof graspnet: Variational grasp generation for object manipulation. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2019)
47. Napier, J.R.: The prehensile movements of the human hand. *The Journal of bone and joint surgery* **38**(4), 902–913 (1956)
48. Oberweger, M., Wohlhart, P., Lepetit, V.: Generalized feedback loop for joint hand-object pose estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* **42**(8), 1898–1912 (2020)
49. Oikonomidis, I., Kyriazis, N., Argyros, A.A.: Full dof tracking of a hand interacting with an object by modeling occlusions and physical constraints. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2011)
50. Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A.A., Tzionas, D., Black, M.J.: Expressive body capture: 3D hands, face, and body from a single image. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)

51. Pham, T., Kyriazis, N., Argyros, A.A., Kheddar, A.: Hand-object contact force estimation from markerless visual tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* **40**(12), 2883–2896 (2018)
52. Pirk, S., Krs, V., Hu, K., Rajasekaran, S.D., Kang, H., Yoshiyasu, Y., Benes, B., Guibas, L.J.: Understanding and exploiting object interaction landscapes. *ACM Transactions on Graphics (TOG)* **36**(3), 31:1–31:14 (2017)
53. Pollard, N.S., Zordan, V.B.: Physically based grasping control from example. In: *Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation* (2005)
54. POSER: 3D rendering and animation software. <https://www.posersoftware.com>
55. Pressure Profile Systems Inc. (PPS). <https://pressureprofile.com>
56. Prokudin, S., Lassner, C., Romero, J.: Efficient learning on point clouds with basis point sets. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
57. Ranjan, A., Hoffmann, D.T., Tzionas, D., Tang, S., Romero, J., Black, M.J.: Learning multi-human optical flow. *International Journal of Computer Vision (IJCV)* (jan 2020)
58. Rogez, G., Supančič III, J.S., Ramanan, D.: Understanding everyday hands in action from RGB-D images. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2015)
59. Romero, J., Kjellström, H., Kragic, D.: Hands in action: real-time 3D reconstruction of hands in interaction with objects. In: *IEEE International Conference on Robotics and Automation (ICRA)* (2010)
60. Romero, J., Tzionas, D., Black, M.J.: Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics (TOG)* **36**(6), 245:1–245:17 (2017)
61. Rosenhahn, B., Schmaltz, C., Brox, T., Weickert, J., Cremers, D., Seidel, H.: Markerless motion capture of man-machine interaction. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2008)
62. Ruff, H.A.: Infants’ manipulative exploration of objects: Effects of age and object characteristics. *Developmental Psychology* **20**(1), 9 (1984)
63. Sahbani, A., El-Khoury, S., Bidaud, P.: An overview of 3D object grasp synthesis algorithms. *Robotics and Autonomous Systems (RAS)* **60**(3), 326336 (2012)
64. Savva, M., Chang, A.X., Hanrahan, P., Fisher, M., Nießner, M.: SceneGrok: Inferring action maps in 3D environments. *ACM Transactions on Graphics (TOG)* **33**(6), 212:1–212:10 (2014)
65. Savva, M., Chang, A.X., Hanrahan, P., Fisher, M., Nießner, M.: PiGraphs: Learning interaction snapshots from observations. *ACM Transactions on Graphics (TOG)* **35**(4), 139:1139:12 (2016)
66. Sridhar, S., Mueller, F., Zollhoefer, M., Casas, D., Oulasvirta, A., Theobalt, C.: Real-time joint tracking of a hand manipulating an object from RGB-D input. In: *Proceedings of the European Conference on Computer Vision (ECCV)* (2016)
67. Starke, S., Zhang, H., Komura, T., Saito, J.: Neural state machine for character-scene interactions. *ACM Transactions on Graphics (TOG)* **38**(6), 209:1–209:14 (2019)
68. Stratasys Fortus 360mc: 3D printing. <https://www.stratasys.com/resources/search/white-papers/fortus-360mc-400mc>
69. Sundaram, S., Kellnhofer, P., Li, Y., Zhu, J.Y., Torralba, A., Matusik, W.: Learning the signatures of the human grasp using a scalable tactile glove. *Nature* **569**(7758), 698–702 (2019)

70. Tekin, B., Bogo, F., Pollefeys, M.: H+O: Unified egocentric recognition of 3D hand-object poses and interactions. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
71. Tekscan grip system: Tactile grip force and pressure sensing. <https://www.tekscan.com/products-solutions/systems/grip-system>
72. Tsoli, A., Argyros, A.A.: Joint 3d tracking of a deformable object in interaction with a hand. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)
73. Tzionas, D., Ballan, L., Srikantha, A., Aponte, P., Pollefeys, M., Gall, J.: Capturing hands in action using discriminative salient points and physics simulation. *International Journal of Computer Vision (IJCV)* **118**(2), 172–193 (2016)
74. Varol, G., Romero, J., Martin, X., Mahmood, N., Black, M.J., Laptev, I., Schmid, C.: Learning from synthetic humans. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
75. Vicon Vantage: Cutting edge, flagship camera with intelligent feedback and resolution. <https://www.vicon.com/hardware/cameras/vantage>
76. Wang, Y., Min, J., Zhang, J., Liu, Y., Xu, F., Dai, Q., Chai, J.: Video-based hand manipulation capture through composite motion control. *ACM Transactions on Graphics (TOG)* **32**(4), 43:1–43:14 (2013)
77. Wang, Z., Chen, L., Rathore, S., Shin, D., Fowlkes, C.: Geometric pose affordance: 3D human pose with scene constraints. arXiv:1905.07718 (2019)
78. Welschhold, T., Dornhege, C., Burgard, W.: Learning manipulation actions from human demonstrations. In: IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (2016)
79. XSENS: Inertial motion capture. <https://www.xsens.com/motion-capture>
80. Yamamoto, M., Yagishita, K.: Scene constraints-aided tracking of human body. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2000)
81. Ye, Y., Liu, C.K.: Synthesis of detailed hand manipulations using contact sampling. *ACM Transactions on Graphics (TOG)* **31**(4), 41:141:10 (2012)
82. Zhang, H., Bo, Z.H., Yong, J.H., Xu, F.: InteractionFusion: real-time reconstruction of hand poses and deformable objects in hand-object interactions. *ACM Transactions on Graphics (TOG)* **38**(4), 48:1–48:11 (2019)
83. Zhang, Y., Hassan, M., Neumann, H., Black, M.J., Tang, S.: Generating 3D people in scenes without people. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2020)

GRAB: A Dataset of Whole-Body Human Grasping of Objects **Supplemental Material**

Omid Taheri, Nima Ghorbani, Michael J. Black, and Dimitrios Tzionas

Max Planck Institute for Intelligent Systems, Tübingen, Germany
{otaheri,nghorbani,black,dtzionas}@tuebingen.mpg.de

<https://grab.is.tue.mpg.de>

S.1 Supplementary Video

The video on our website presents:

- a narrated overview of our method,
- a wide variety of GRAB sequences (3D moving meshes),
- GrabNet predictions for unseen objects from several viewpoints and
- GrabNet failure cases.

S.2 GRAB Dataset Content

Our GRAB dataset is available for research purposes on our website. The website contains (at least) the material listed below:

- Our modified version of the object meshes of [2].
- Our marker locations on each object mesh.
- Body shape templates for our subjects.
- Pose parameters for our subjects and objects.
- Code to reproduce the interacting meshes, as seen in our video.
- Per-vertex contact annotations on meshes (body and object) for each frame.
- Vicon MoCap files (labeled marker positions, incl. on the floor and table).

S.3 Why MoCap Instead of 3D Scan Sequences?

For accurate human shapes, we capture a dense 3D scan for each subject to which we fit a personalized 3D SMPL-X template mesh. However, 3D scanning does not scale up for capturing human-object interaction sequences. This would produce huge amounts of data, the processing of which would be a major undertaking. Moreover, object tracking under occlusions would be still very challenging, as finding scan-to-model correspondences is a hard ill-posed problem. Instead, with MoCap a minimum of 3 marker observations is enough for reliable object pose estimation. The placement of many small markers on the objects means that we

can always estimate object pose. Using MoSh++ for the body, given a ground truth body shape, produces accurate meshes that are on par with 3D scanning but much more practical to capture. We follow therefore this practical and scalable approach; we use a high-end optical MoCap system (Sec. 3.1) and fit full 3D meshes to MoCap markers (Sec. 3.2) for both the human and the object.

S.4 Why not Capture RGB Images?

Capturing *accurate* human-object interactions while also capturing *natural* RGB images is very challenging. Some recent datasets [4, 5] capture hand-only interactions with objects and include RGB images, but the images capture only the hand and not the whole body [4, 5] and are not fully natural due to visible instrumentation on the hand [4]. Please note that this latter point is fundamental. Currently one must choose between accurate grasping, which requires instrumentation, or natural images, which reduces the accuracy of ground truth.

Both methods [4, 5] suffer from severe hand-object inter-penetrations. Garcia-Hernando et al. [4] originally reconstruct a hand skeleton interacting with 4 object meshes, and their method was reported to have an average *skeleton* penetration depth of 11.0 ± 8.9 mm (see Sec. 5.2 of [6]). Similarly, we compute the *surface* penetration between MANO and the 3D object meshes for [5] and find the mean to be 4.36 ± 0.94 mm. Although the hand inter-penetration of [5] is not as severe as [4], it suffers from not having realistic contact with objects. In Fig. S.1 we compare the contact “heatmaps” for [5] (left) and GRAB (right). Note that for [5], the thumb and all fingertips are rarely in contact, whereas for GRAB they are frequently in contact. The latter is much more realistic given the central role of the thumb and fingertips in object manipulation. This points to an important technical problem, without lowering the value of these works, as they focus on a challenging application.

We conclude that state-of-the-art interaction methods, that also capture RGB images, suffer from intense occlusions and penetrations along with non-realistic contact between the hand and the objects. Such data is not good to learn an accurate data-driven model of 3D interactions. In contrast, our “use” grasps have only 3.25 ± 0.68 mm average *surface* penetration, which is significantly lower than [4, 5], and realistic contact between the body and objects, while containing more challenging scenarios, namely dexterous in-hand manipulation and capture of the whole body instead of only the hand.

This is attributed to our precision-focused setup, that increases accuracy on the expense of not capturing RGB images, due to the uniform and artificial texture of the MoCap body suit and the 3D printed objects. We believe that this is a sensible trade-off; one can use our accurate 3D mesh reconstructions to learn a model of 3D interactions, and use it as a prior in future work to improve methods like [4, 5] for the hand or the whole body.

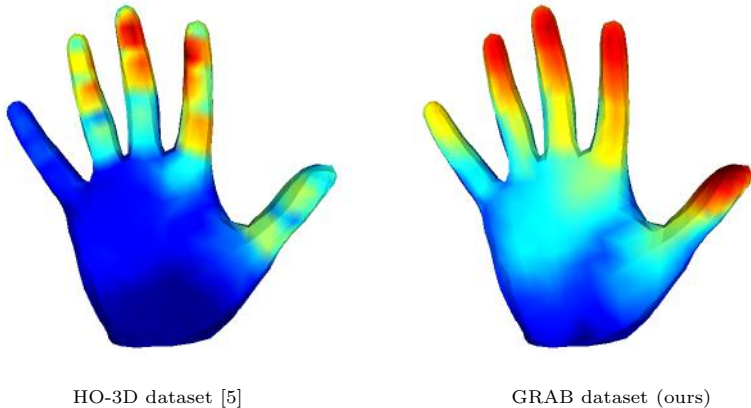


Fig. S.1: Contact “heatmaps” for HO-3D [5] (left) and GRAB (right), for the right hand. The hotter the color, the more frequently that hand part is in contact with objects. During grasping and manipulation, the thumb finger and all fingertips play a central role. This is evident with GRAB but not with [5].

S.5 Penetration Plots

In the main paper we describe the observed human-object mesh inter-penetration. Due to space constraints, we report here the corresponding plots in Fig. S.2. We evaluate the degree of penetration for “use” sequences, that pose the most realistic occlusions and capture challenges. “Use” grasps have 3.25 ± 0.68 mm average penetration, which effectively corresponds to the missing soft-tissue deformation. Please note that there is no model of the human body with articulated fingers and face that captures such soft-tissue deformation with contact. In addition, 67% of “use” grasps have ≤ 3.5 mm penetration, 86% ≤ 4.0 mm, 96% ≤ 4.5 mm and 99.9% has ≤ 5.8 mm.

S.6 Protocol Details

Here we provide details that were not crucial for the main manuscript (Sec 3.4). We capture motions with 4 intents: “use”, “pass”, “lift”, and “off-hand pass”. For each sequence we randomize the object position and pose on a resting table, the height of which is also randomized between 75 cm and 120 cm to increase motion variance. We capture the following intents:

“Use”: For the objects that have a clear everyday use (e.g. drinking from cup), we ask the subject to naturally use them. In case of multiple use cases (e.g. digital/analog photo camera) we capture multiple sequences. For objects without a clear use (e.g. cylinder) the subject has to grasp them and inspect them.

“Pass”: The subject is asked to pass the object to a predefined direction, that is randomized (e.g. bottom-left, top-right, etc), to increase motion variance.

“Lift”: The subject is asked to grasp the object, lift it stably in any natural way

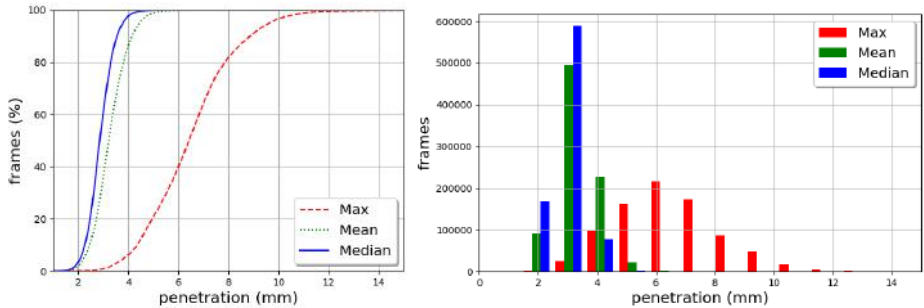


Fig. S.2: Penetration plots for “use” grasps. For each frame we store the max (red), mean (green) and median (blue) vertex penetration. **(Left):** percentage of frames (Y axis) below a varying penetration error (X axis). **(Right):** bar-plot for number of frames (Y axis) with a specific (quantized for binning) penetration (X axis). The mean penetration is 3.25 ± 0.68 mm.

they can imagine, then leave it on the table in any random pose, and repeated this several times with re-grasping. This increases grasp variance, by encouraging the exploration of contact configurations and relative hand-object orientations.

“Off-Hand pass”: As a form of bimanual manipulation, the subject grasps the object with the off hand, passes it to the dominant hand, and uses it (see “use”).

We capture MoCap markers placed on the body, face and fingers, as well as on the object (Sec 3.1 in paper). Additionally, we capture markers attached to the floor and the table, for potential future use. All subjects gave informed consent to share their motion data for research purposes.

S.7 Computing Contact

Here we provide some additional intuition to Sec. 3.3 and Fig. 4 (right). In particular we explain how we deal with noise in the reconstructed moving meshes to produce clean contact data.

Figure S.3 (left): Consider the illustrated example of a 3D cup. Its mesh thickness is thin, i.e. it has an outer and inner surface that are different, yet close to each other. In Fig. S.3 (top-left) the thumb and index fingers of a grasping hand penetrate both the outer and inner surface. This is due to noise, fitting errors, and because existing models do not model contact-dependent skin deformations.

For these examples the actual contact area is the one on the outer object surface. To annotate only this, following Sec. 3.3, we first compute all colliding triangles and cluster them in connected “rings” (Fig. S.3, top-left). For each “ring” we compute the corresponding penetrating hand areas (Fig. S.3, bottom). The hand areas that contact the inner surface are a subset of the ones that contact the outer one. Then, we remove (big red circles in Fig. S.3 bottom) the

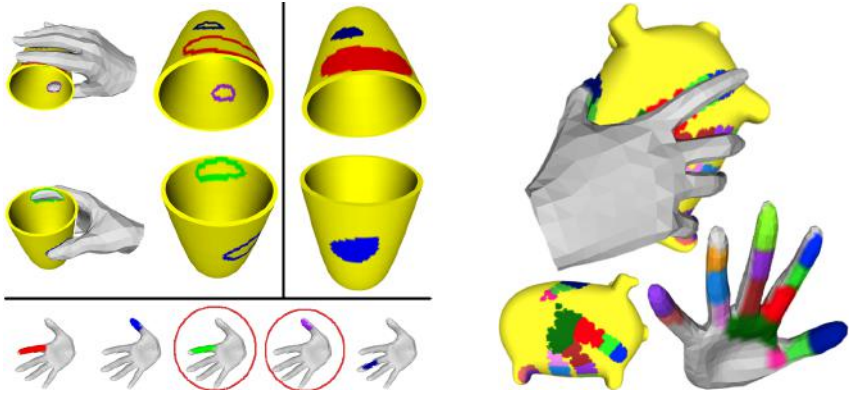


Fig. S.3: (Left): Annotating contact areas for a hand grasping a cup: (top-left) “rings” of colliding triangles, color-coded for each finger, (bottom) penetrating hand areas that correspond to each “ring”; the ones corresponding to the green and purple vertices (red circled hand parts) penetrate the inner cup surface and are ignored, (top-right) the final filtered “rings” and the enclosed vertices are annotated as contact areas. The contact labels are binary; color is used here only for visualization purposes. **(Right):** Contact labels can though be more fine-grained, e.g. using the contacting hand parts or hand vertices. Here we see an example of the former case. (top) Each color represents a contact area caused by a different hand part. (bottom) Contact areas are shown also on the object and unposed hand for clarity. Note that the size of contact areas is expanded for illustration purposes.

purple and green groups, we keep only the remaining “rings”, and annotate the vertices enclosed by them as contact vertices (Fig. S.3, top-right).

Figure S.3 (right): The above procedure gives binary contact annotations (“contact” on “not in contact”). Contact labels, however, can be more fine-grained, e.g. with the label of the corresponding hand part (Fig. S.3, right), or even with the point on the 3D hand surface (Fig. S.4). For the former example, we find the object vertices that are in contact, and for each one we find the closest SMPL-X/MANO bone, and assign its ID as the contact label.

S.8 Adapting MoSh++

We adapt MoSh++ [9] (Sec. 3.2) for capturing the whole body (including the hands and face). The human and object are tracked independently and on a per-frame basis, for simplicity. We make two small changes to MoSh++. First we use the ground-truth body shape, obtained from a 3D scan. Consequently, we do not use MoSh++ to estimate body shape. Second, we extend MoSh++ to estimate the parameters of the SMPL-X body model. This means extending it to capture facial pose and expression parameters. Additionally, we estimate the rigid 6 DoF pose of the objects using their known shape and the detected markers.

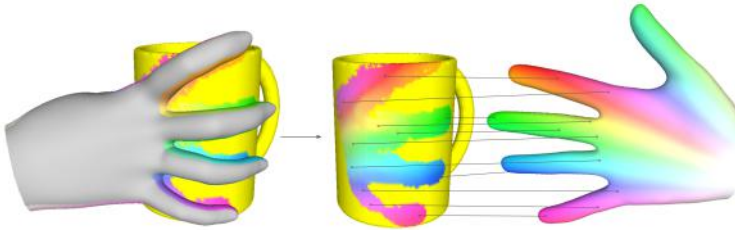


Fig. S.4: (Right): Fine-grained contact labels. In contrast to the binary contact labels of Fig. 4 (left) and Fig. 9 of the main manuscript, and the part-based contact labels of Fig. S.3, here we show an example of much more fine-grained labels. **(left)** A 3D hand-object grasp configuration. **(middle)** The object alone. **(right)** The hand in canonical pose. We highlight different points on the 3D surface of the inner hand with color gradients. The contact between the hand and object define surface correspondences between them (shown as lines).

To adapt MoSh++ to capture faces, we need to tune the parameters of the model. For this [9] follows a data-driven approach; they capture the SSM dataset with an optical MoCap system synchronized with a 3D body scanner, and use the scans for computing a reconstruction quality metric. However, SSM has markers only on the main body, while also the fingers of the scans are very noisy.

Capturing such a dataset, with clear scan regions for both the body, the face and all fingers, as well as synced MoCap for them, is too challenging. Instead, we follow a more practical approach, and create a synthetic dataset by animating SMPL-X and generating virtual markers on the moving meshes. To bridge the domain gap, we simulate noise for marker position and visibility; we randomly add 3D Gaussian noise with 1 mm variance in marker positions, as in [9], and randomly drop up to 5 markers per frame.

Unfortunately, there is no existing dataset with rich SMPL-X sequences. However, its model formulation is compatible to SMPL [8] for the body, FLAME [7] for the face, and MANO [13] for the hands. Therefore, we resort to datasets specific to each part to animate the body, face and hands. For the *body*, we employ DFAUST [1] that captures 10 subjects performing 10 sequences each. We split the subjects into 6 for training and 4 for a withheld test set. We compute personalized SMPL-X mesh templates by registering the model to one scan per person as in [13, 10], and pose their body according to the registrations of DFAUST. For the *hand*, we employ the hand-only MANO model registrations of [13]. From the 1554 hand poses we hold out 155 for the test set and use the rest for training. We then add hand motion to each body sequence by randomly choosing 15 hand poses and interpolating between them. For the *face*, we employ sequences of FLAME parameters from [3, 12]; the latter covers extreme facial expressions, while the former has everyday speaking expressions. We randomly choose 100 sequences from each dataset, splitting them into 60 for the training set and 40 for the withheld test set.

Table S.1: Evaluation of MoSh++ on the synthetic dataset. We compare the vanilla [9] to our adapted version. For the first stage of MoSh++, Stage-I, we report the distance of the latent marker placement compared to ground-truth marker locations, and for the second stage of MoSh++, Stage-II, we report the average vertex-to-vertex error between estimated and ground-truth meshes.

Mosh++ version	MoSh Stage-I		MoSh Stage-II	
	mean \pm std	median	mean \pm std	median
Vanila [9]	4.76 \pm 1.03	4.55	5.59 \pm 1.86	5.28
Our adapted	3.09 \pm 0.55	2.80	4.86 \pm 1.83	4.48

mm

We use this dataset to set the weights following the approach in [9]. Table S.1 compares a standard version of [9] with our adapted version on the synthetic test set for both stages of MoSh++. For the first stage of MoSh++ (Stage-I) we report the distance of the latent marker placement compared to ground-truth marker locations in mm. For this stage we start from random marker placement guesses in the 1-ring neighborhood of the ground-truth locations. We repeat this three times with different random seeds for selecting 12 frames of MoSh++; see [9]. For the second stage (Stage-II), we use the optimized latent marker placements resulting from each random seed of the Stage-I and report the average vertex-to-vertex error between estimated and ground-truth meshes in mm. In each stage we choose the wights that minimizes the reported error. Our adapted version shows a clear improvement, by fitting the whole body, hands and face, with weights λ tuned on our synthetic dataset. In contrast [9] tunes only the body weights on their SSM dataset, it fits the hands with empirical weights, and does not fit the face.

S.9 GrabNet

To show that GRAB can be useful for machine learning applications, we train on it a generative model to generate grasping poses for a 3D object, that we call GrabNet. For this example application we focus only on right-hand grasps for simplicity, but GRAB provides much richer data than this.

S.9.1 Data Preparation

We use only right hand data due to the large size of GRAB, but left hand poses could also be mirrored to appear as right ones for data augmentation. In order to select right-hand frames for training GrabNet we use the following rules. (i) The right hand should be in contact. (ii) The left hand should not have any contact. (iii) The object’s vertical position should be at least 5 mm different from its initial one (i.e. it should be lifted from the resting table). (iv) The right thumb and at least one more finger should be in contact. (v) A finger is considered a contacting finger, when it is in contact with at least 50 object vertices. With

these filters we make sure that we have only stable grasps with which to train GrabNet.

To model arbitrary shapes, we use the basis point set [11] representation BPS_o for all our objects. For computational efficiency, we precompute BPS_o and load it from memory during training. We sample basis points in a sphere of 150 mm radius, that is big enough to cover our centroid-centered objects. We empirically found 4096 basis points to be enough. We then compute the distances from the basis points to our object meshes.

Out of our 51 objects, borrowed from [2], we hold out 4 objects for the validation set (“apple”, “toothbrush”, “elephant” and “hand”), 6 objects for the test set (“mug”, “wineglass”, “camera”, “binoculars”, “frying pan” and “toothpaste”), and use the remaining 41 objects for the training set.

The training, validation and test splits contain roughly 320k, 31k and 65k data points, correspondingly.

To prepare the training data for RefineNet we add Gaussian noise to the Ground Truth MANO parameters of the selected data for GrabNet. Since the perturbation need to be minimal we empirically find $\mathcal{N}(\mu = 0, \sigma^2 = 0.2)$, $\mathcal{N}(\mu = 0, \sigma^2 = 0.004)$, and $\mathcal{N}(\mu = 0, \sigma^2 = 0.05)$ for MANO finger joints rotation, root rotation, and translation respectively.

S.9.2 Results: Success and Failure Cases

Figures S.6, S.7 and S.8 provide a wide variety of qualitative GrabNet results. More specifically, they show 10 different grasps (rows) generated for 6 unseen objects (columns). The three figures show three different viewpoints (one view per figure) for the same grasp of the 10×6 grid. We see that most grasps look natural and plausible, as GrabNet is learned from high-quality GRAB captures. More results with a rotating viewpoint are shown in the video on our website.

GrabNet can still generate some failure cases. These are mostly cases of penetrating fingers; there are not many cases of contacting fingers that fly away from the object. Penetrations are observed mostly for objects with thin parts (cup handle, wine glass, bowl). We found the frying pan to be the most challenging object, due to its comparably big size along with its thin surface walls and handle. This might be due to the sparse BPS_o representation for 3D object shapes capturing mostly their bigger parts. Furthermore, at the moment we use a penetration and a contact term in the training loss of GrabNet as soft constraints, since here we focus on a data-driven method. One could add an optimization stage to refine the regressed grasp with hard penetration and contact constraints.

The results show the value of GRAB for training data-driven models, but also point to room for improvement for GrabNet’s modeling and training scheme.

S.9.3 GrabNet Implementation Details

The architecture for GrabNet is shown in Fig. S.5. For CoarseNet, we concatenate the object BPS_o representation with MANO hand parameters as input to the

encoder, and also concatenate it with the latent code as the input (condition) to the decoder. The outputs of the decoder are MANO translation ($\gamma \in R^3$) and joint angles ($\theta \in R^{96}$) in the continuous 6-dimensional representation of [15].

Using our validation set, we found out that 16 dimensions for the latent space results in generating better grasps. Qualitative results are provided in Fig S.6.

For RefineNet we take the output of the CoarseNet (MANO parameters) and first compute the distances of MANO vertices to the object vertices. We then pass the distances with the MANO parameters to the network. RefineNet refines the input grasp through 3 iterations. The CoarseNet and RefineNet are trained for 16 and 23 epochs respectively with the learning rate starting from $5e - 4$, decreasing on validation error plateau to 0.1 times, and early stopping after 8 epochs with no improvement in validation error. Both networks are trained separately.

S.9.4 Filtering out Unreliable Turkers

As mentioned in the main paper (Sec. 4), along with ground-truth (GRAB) and GrabNet-generated grasps, we pass to Turkers noisy grasps generated by perturbing ground-truth ones. These noisy grasps are our test for spotting unreliable Turkers, that either select their answers randomly or misunderstand the task. Specifically, we remove the ones that gave a rating of 3 or more (indicating good realism) for at least 20% of these noisy grasps. In total we removed 54 out of 170 Turkers.

S.9.5 Heatmaps for Various Intents and Fine-Grained Numbers

Similar to Sec 3.5 of the main manuscript, here Fig. S.9 provides additional fine-grained numbers for in-contact parts of the body. Each row corresponds to an intent in the GRAB dataset. For each intent, the right column shows the contact percentage and “heatmap” for the right hand, left hand, and head across all frames and relative to *all* body vertices. In the three left columns, the “heatmaps” and percentages are relative to only *each part’s* vertices and for only the frames for which these parts are in contact (left hand, right hand, and head), for visualization purposes. For example, for the “use” sequences (second row in Fig S.9), the right hand was in contact for 90.62% of all frames, and in those contact frames the thumb fingertip was in contact for 99.88% of them.

S.10 Bias from MoCap Markers

A natural question arises - are subjects biased in their grasps by MoCap markers? We empirically place more markers in areas less likely to be contacted, according to object affordances. To account for potential occlusions, though, we have to place some markers in other areas as well. For this reason, we still expect our markers to be contacted.

Our subjects did not complain about discomfort or bias, yet we need more evidence for this. Apart from the analysis in the main manuscript (see Sec 3.5 and Fig. 5), here we perform k-means clustering ($k=20$) on our grasps, and visualize each cluster center, i.e. a grasping hand, and the grasped object. We observe that several clusters (typically 3-6 out of 20) show that fingers do come in contact with markers. Figure S.10 shows for 5 objects (rows) 3 representative contacting clusters (columns). We believe this is good additional empirical evidence that our 1.5 mm radius hemi-spherical markers cause no or minimal bias.

S.11 Influence of Contact Heuristic Thresholds

We use several heuristics to determine contact frames, see Sec. 3.3 in the main manuscript. For the contact “heatmap” analysis we take all the contact frames for which object is being manipulated, i.e. it is off the table. Because the heatmap is integrated over many frames, small variations in the heuristics have little impact on the contact patterns.

To show this empirically, we perform a sensitivity analysis by changing our thresholds. Figure S.11 shows “heatmaps” for “use” sequences for several setups, following the format of Fig. 5 of the main manuscript. The results verify our hypothesis that the heuristics have minimal influence.

S.12 Influence of Object Size

Figure S.12 shows the effect of object size on hand contact for 5 shape primitive objects. Smaller objects are grasped mainly by the dominant hand, while, for larger objects, the other hand helps too. Small-sized subjects have a large tendency for bimanual grasps, and involve more finger and contact around metacarpophalangeal joints.

S.13 Vicon Software

We use the Vicon system with the Shōgun [14] software, in two main modes. Shōgun-Live is used for camera calibration and synchronization, as well as for subject calibration and data capture. The output consists of 3D marker locations and labels (IDs). However, both can be noisy; markers can disappear due to extreme occlusions, ghost markers might appear due to reflections, while labels can be mistakenly swapped when many markers come close together. This calls for the need for data cleaning in a post-processing stage. This is important for strong interactions, as getting even only 3 cleaned marker observations is enough to solve for the rigid pose of a manipulated object. Shōgun-Post is used for such semi-automatic data cleaning by 4 trained annotators.

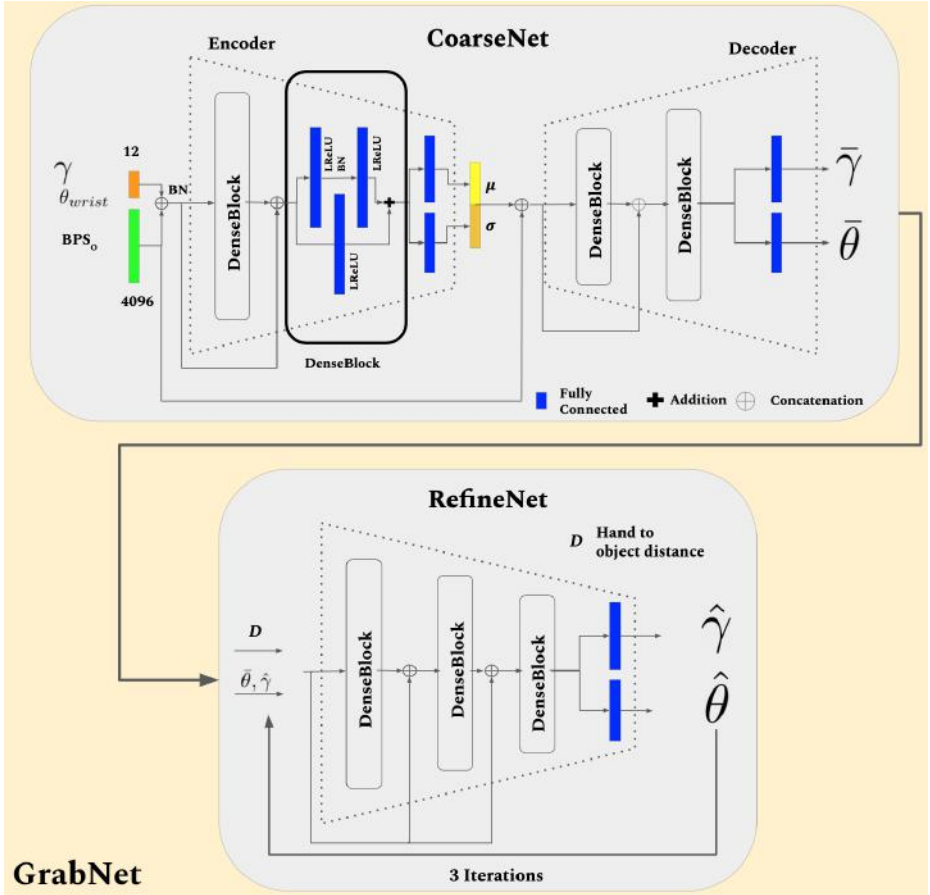


Fig. S.5: GrabNet Architecture. For the encoder input, we concatenate the BPS representation of the object with MANO parameters, while for decoder input we concatenate it with a sample from latent space. The decoder gives the MANO hand parameters which we pass to the MANO model to obtain the 3D hand mesh.

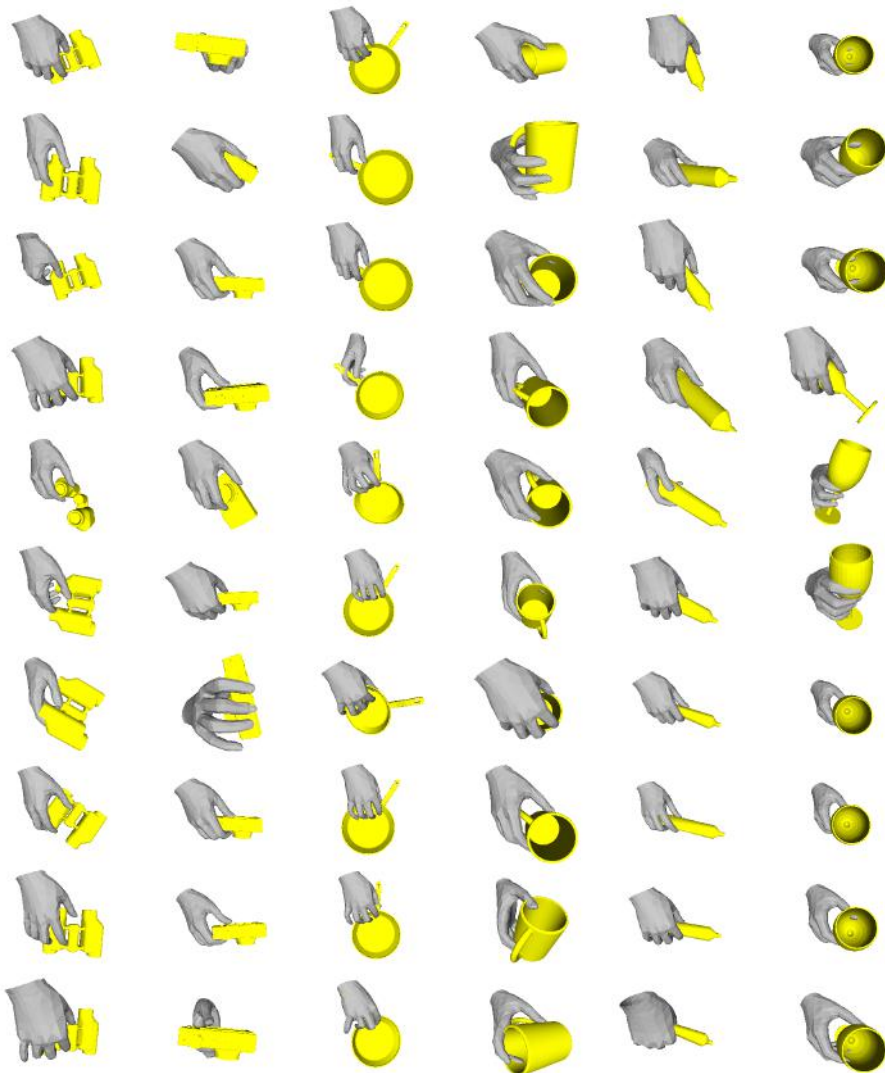


Fig. S.6: Visualization of 10 different grasps (rows) generated by GrabNet for 6 unseen objects (columns). Conditioned on the BPS_o representation of unseen 3D object shapes, we sample from the learned 16 dimensional latent grasping space Z of CoarseNet (see Fig. 8 of the paper). We then concatenate BPS_o to the Z sample, and pass them to the decoder of CoarseNet, that outputs the coarse grasping MANO hand model parameters. We then pass the coarse grasps to the RefineNet to get the final grasps. Some failure cases are highlighted with red. Different viewpoint for the results of Fig. S.7, S.8.

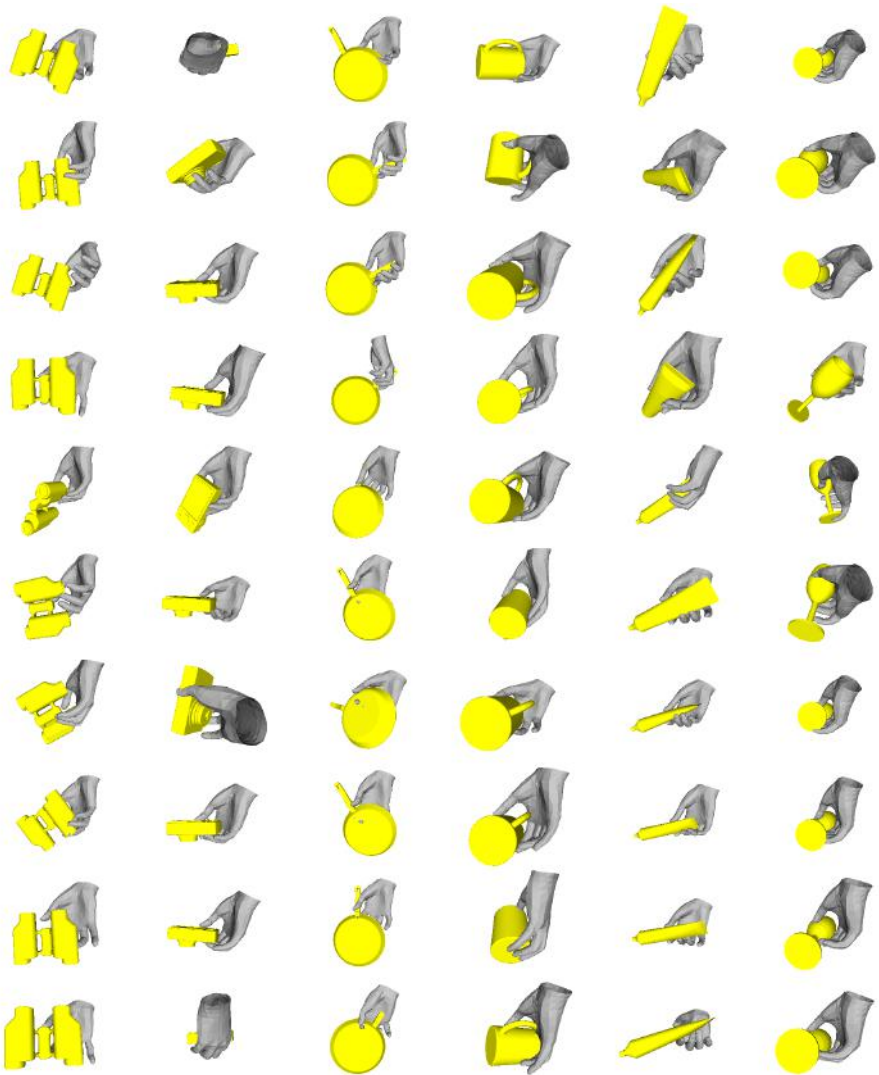


Fig. S.7: Visualization of 10 different grasps (rows) generated by GrabNet for 6 unseen objects (columns). Conditioned on the BPS_o representation of unseen 3D object shapes, we sample from the learned 16 dimensional latent grasping space Z of CoarseNet (see Fig. 8 of the paper). We then concatenate BPS_o to the Z sample, and pass them to the decoder of CoarseNet, that outputs the coarse grasping MANO hand model parameters. We then pass the coarse grasps to the RefineNet to get the final grasps. Some failure cases are highlighted with red. Different viewpoint for the results of Fig. S.6, S.8.

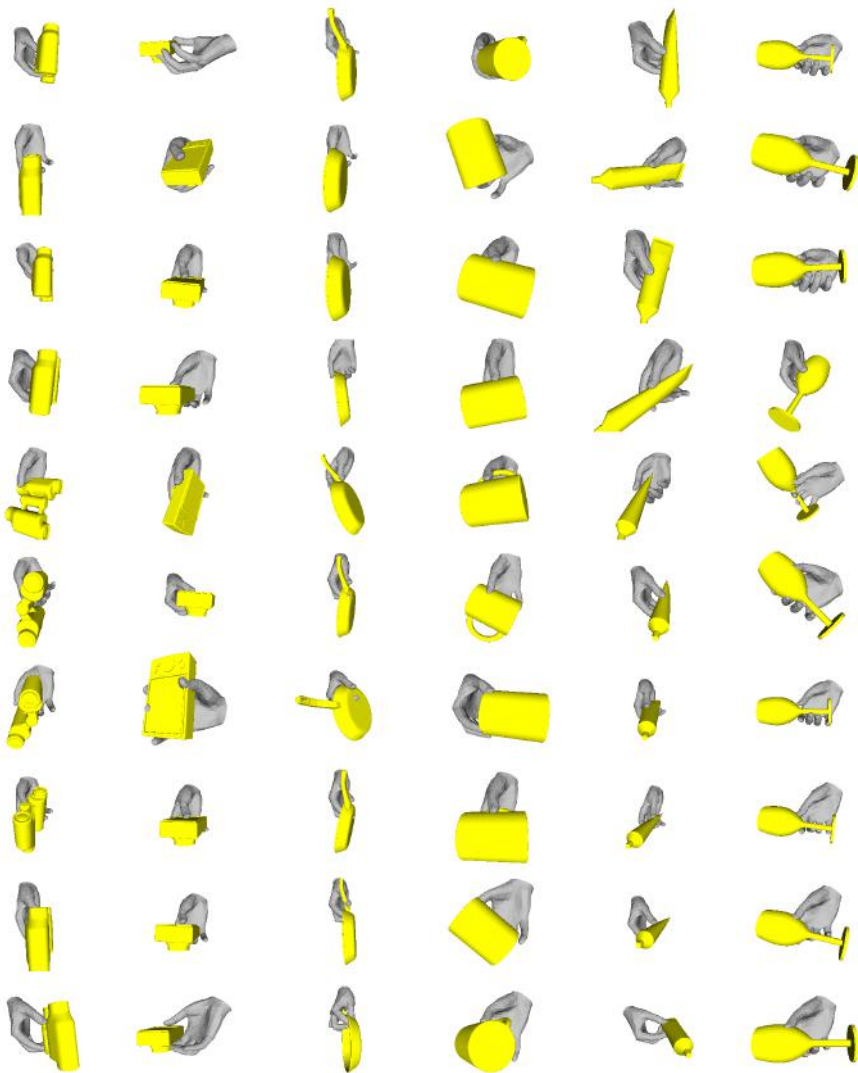


Fig. S.8: Visualization of 10 different grasps (rows) generated by GrabNet for 6 unseen objects (columns). Conditioned on the BPS_o representation of unseen 3D object shapes, we sample from the learned 16 dimensional latent grasping space Z of CoarseNet (see Fig. 8 of the paper). We then concatenate BPS_o to the Z sample, and pass them to the decoder of CoarseNet, that outputs the coarse grasping MANO hand model parameters. We then pass the coarse grasps to the RefineNet to get the final grasps. Some failure cases are highlighted with red. Different viewpoint for the results of Fig. S.6, S.7.

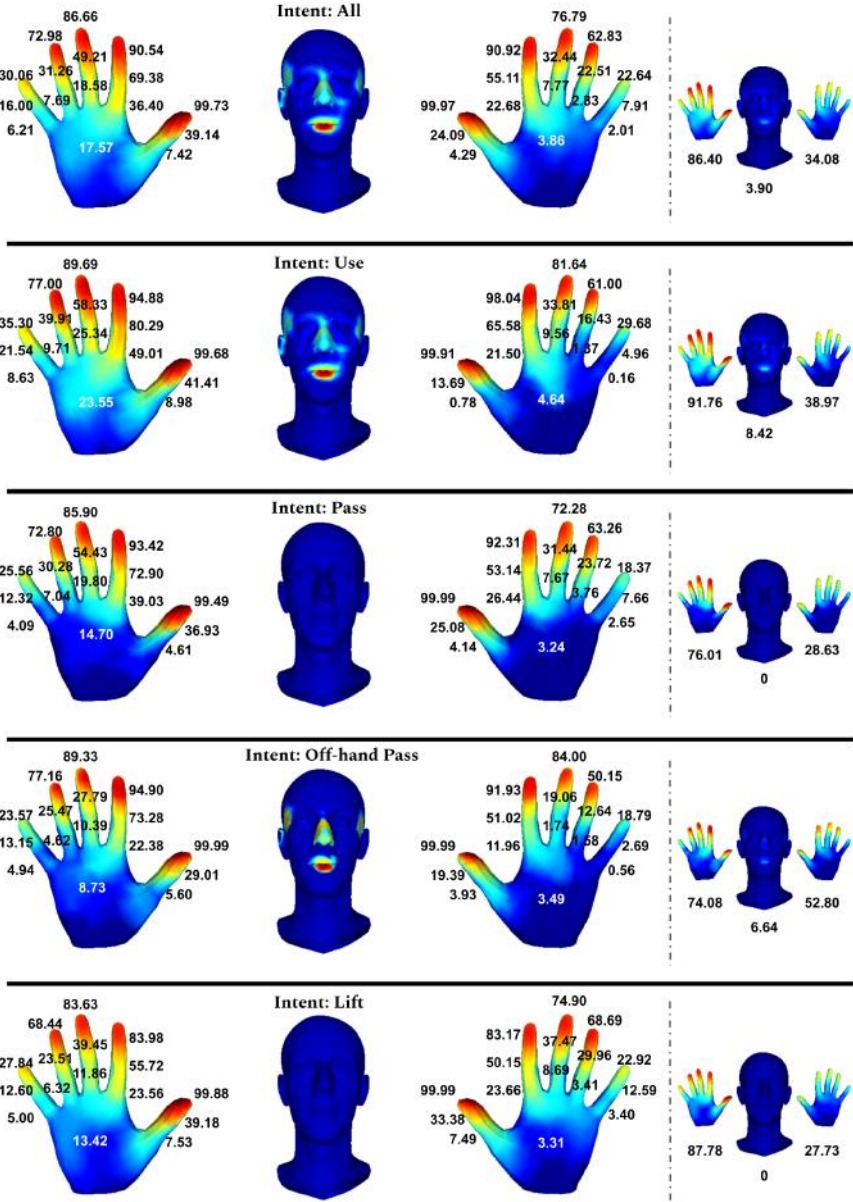


Fig. S.9: Contact “heatmaps” and percentages for all intents in GRAB, for various body parts. Each row corresponds to an intent in the GRAB dataset. For each intent, the right column shows the results for each part (right hand, left hand, and head) across all frames and relative to *all* body vertices. In the three left columns, the results are relative to only *each part’s* vertices and for only frames for which these parts are in contact (left hand, right hand, and head), for visualization purposes.

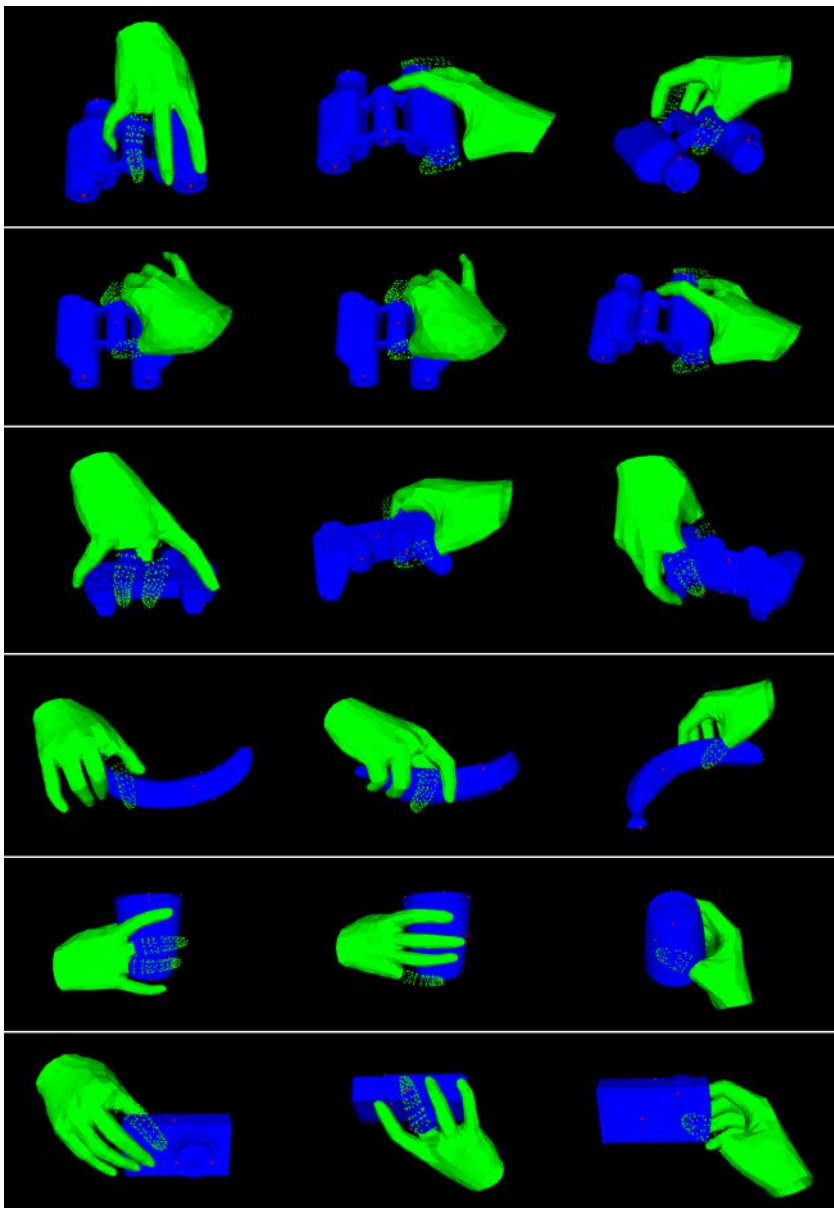


Fig. S.10: Do subjects avoid markers? To answer this, we perform k-means clustering ($k=20$) on our grasps, and visualize each cluster center, i.e. a grasping MANO (green), and the grasped object (blue). We observe that several cluster centers (columns) per object (rows) show that subjects contact MoCap markers (red); here we show 3 clusters for 5 objects. For fingers that contact markers we render only the vertices, to allow to see the markers (best viewed on screen).

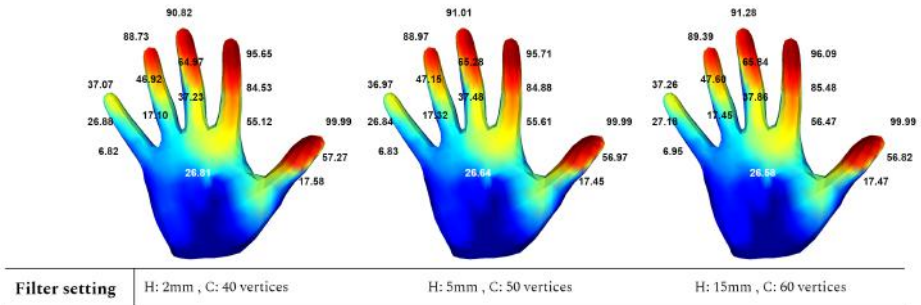


Fig. S.11: Sensitivity analysis for contact heuristics. We follow the format of Fig. 5 of the main manuscript, and show “heatmaps” and contact likelihoods in percentages % for a subset of “all” sequences with different setups (columns), as indicated in the labels. The symbol H denotes the minimum difference between the object’s vertical position from its initial one (resting on a table). The symbol C denotes the minimum number of object vertices that we require to be in contact with each finger. The figure shows that threshold choices have minimal effect when integrated over many frames to create “heatmaps”.

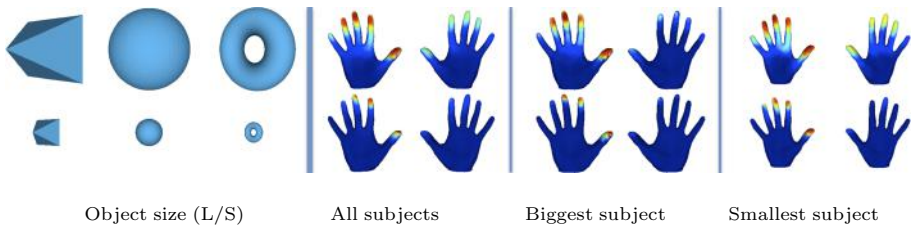


Fig. S.12: Effect of object size on contact during grasping. We show contact “heatmaps” for all subjects, the biggest subject, and the smallest subject.

References

1. Bogo, F., Romero, J., Pons-Moll, G., Black, M.J.: Dynamic faust: Registering human bodies in motion. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
2. Brahmabhatt, S., Ham, C., Kemp, C.C., Hays, J.: ContactDB: Analyzing and predicting grasp contact via thermal imaging. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
3. Cudeiro, D., Bolkart, T., Laidlaw, C., Ranjan, A., Black, M.: Capture, learning, and synthesis of 3D speaking styles. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
4. Garcia-Hernando, G., Yuan, S., Baek, S., Kim, T.K.: First-person hand action benchmark with rgb-d videos and 3d hand pose annotations. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
5. Hampali, S., Oberweger, M., Rad, M., Lepetit, V.: HO-3D: A multi-user, multi-object dataset for joint 3D hand-object pose estimation. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
6. Hasson, Y., Varol, G., Tzionas, D., Kalevatykh, I., Black, M.J., Laptev, I., Schmid, C.: Learning joint reconstruction of hands and manipulated objects. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
7. Li, T., Bolkart, T., Black, M.J., Li, H., Romero, J.: Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics (TOG)* **36**(6), 194:1–194:17 (2017)
8. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. *ACM Transactions on Graphics (TOG)* **34**(6), 248:1–248:16 (2015)
9. Mahmood, N., Ghorbani, N., F. Troje, N., Pons-Moll, G., Black, M.J.: AMASS: Archive of motion capture as surface shapes. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2019)
10. Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A.A., Tzionas, D., Black, M.J.: Expressive body capture: 3D hands, face, and body from a single image. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
11. Prokudin, S., Lassner, C., Romero, J.: Efficient learning on point clouds with basis point sets. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
12. Ranjan, A., Bolkart, T., Sanyal, S., Black, M.J.: Generating 3D faces using convolutional mesh autoencoders. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)
13. Romero, J., Tzionas, D., Black, M.J.: Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics (TOG)* **36**(6), 245:1–245:17 (2017)
14. Vicon Shogun: VFX motion capture. <https://www.vicon.com/software/shogun>
15. Zhou, Y., Barnes, C., Lu, J., Yang, J., Li, H.: On the continuity of rotation representations in neural networks. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019)