

# EgoPressure: A Dataset for Hand Pressure and Pose Estimation in Egocentric Vision

Yiming Zhao<sup>1\*</sup>

Taein Kwon<sup>1\*</sup>

Paul Strel<sup>1\*</sup>

Marc Pollefeys<sup>1,2</sup>

Christian Holz<sup>1</sup>

<sup>1</sup>ETH Zürich

<sup>2</sup>Microsoft Spatial AI Lab, Zürich

[yiming-zhao.github.io/EgoPressure](https://yiming-zhao.github.io/EgoPressure)

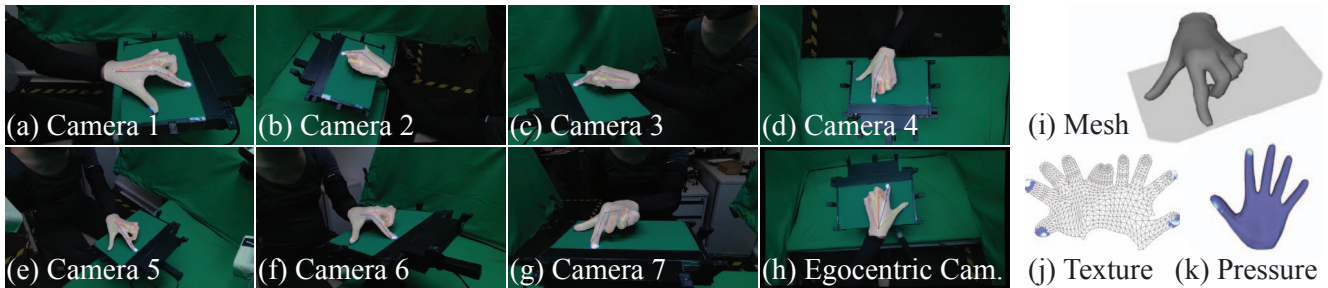


Figure 1. **The EgoPressure dataset.** We introduce a novel egocentric pressure dataset with hand poses. We label hand poses using our proposed optimization method across all static camera views (Cameras 1–7). The annotated hand mesh aligns well with the egocentric camera’s view, indicating the high fidelity of our annotations. We project the pressure intensity and annotated hand mesh (Fig. *i*) to all camera views (Fig. *a* to *h*), and further provide the pressure applied over the hand as a UV texture map (Fig. *j* and *k*).

## Abstract

*Touch contact and pressure are essential for understanding how humans interact with objects and offer insights that benefit applications in mixed reality and robotics. Estimating these interactions from an egocentric camera perspective is challenging, largely due to the lack of comprehensive datasets that provide both hand poses and pressure annotations. In this paper, we present EgoPressure, an egocentric dataset that is annotated with high-resolution pressure intensities at contact points and precise hand pose meshes, obtained via our multi-view, sequence-based optimization method. We introduce baseline models for estimating applied pressure on external surfaces from RGB images, both with and without hand pose information, as well as a joint model for predicting hand pose and the pressure distribution across the hand mesh. Our experiments show that pressure and hand pose complement each other in understanding hand-object interactions.*

## 1. Introduction

Understanding touch during hand-object interaction, especially from an egocentric perspective, is key for augmented reality (AR) [31, 72], virtual reality (VR) [20, 66], and

robotic manipulation [9, 10, 48]. In AR/VR environments, touch contact and pressure information allow for more precise control and feedback [8]. For example, a virtual piano could vary its sound with key pressure, a feature lacking in current AR/VR systems [50]. Pressure sensing is also crucial for robots to replicate human grasping, where precise force estimation remains a challenge [9, 10, 39].

Previous approaches have used gloves [46, 47] and robots with tactile sensors [39, 80] to capture pressure measurements during object manipulation. However, this instrumentation interferes with natural touch by obstructing tactile feedback. In contrast, vision-based estimation methods require no instrumentation of the hands, and cameras are already integrated into devices like smart glasses and mixed reality headsets [21, 22]. Despite this potential, advancements in state-of-the-art models have been limited by the lack of relevant datasets with contact pressure annotations. A notable exception is the PressureVision dataset [19] that comprises RGB footage from four static cameras of hands interacting with a pressure-sensitive surface and corresponding projected pressure images.

In this paper, we introduce a novel dataset, EgoPressure, that extends these prior efforts [19, 20] and captures hand-surface interactions from an *egocentric* perspective, complete with pressure maps projected onto the articulated *hand mesh* in 3-space. Our capture platform combines a Sensel Morph touchpad with one head-mounted and seven static

\* Equal contribution.

cameras, all recording RGB-D data at 30 Hz (Figure 1). The dataset includes 5 hours of footage from 21 participants, each performing 64 interaction sequences with an average length of 420 frames—making it the first bare-handed egocentric dataset with pressure and hand mesh annotations.

We further provide baseline models to demonstrate the potential of our dataset and establish a benchmark for future research. First, we set PressureVisionNet [19] as a baseline on our egocentric dataset and compare it to adapted models that incorporate hand pose as additional input. The model using hand poses estimated from the RGB images via the HaMeR [54] estimator outperforms PressureVisionNet by more than 5% in volumetric IoU error, with improvements of over 7% when using ground-truth hand poses. Additionally, we introduce the first model to jointly estimate hand pose, hand mesh, and pressure both over the mesh and on the surface from an egocentric RGB camera, thereby localizing contact and pressure in 3D space.

We summarize our key contributions as follows:

1. EgoPressure is the first egocentric hand-surface interaction dataset with projection-based pressure annotations together with 3D hand meshes.
2. We establish two novel benchmarks: (1) estimating contact pressure from egocentric RGB images with and without hand pose information, and (2) jointly predicting 3D hand poses and applied pressure, including the localization of pressure on a user’s hand mesh.

EgoPressure thus offers new opportunities for future work to address the unique challenges of egocentric camera views and to precisely localize pressure on a user’s hand.

## 2. Related Work

**Vision-based hand-object pose estimation** Over the past decade, significant progress has been made in hand tracking due to advancements in deep learning techniques [26, 53, 54] and the collection of relevant datasets [52, 76, 79]. While egocentric hand tracking for gesture recognition and direct input has advanced to the point of integration into modern commercial devices such as AR and VR headsets [27, 28], understanding hand interactions with external objects remains an active area of research [15, 16, 22, 38, 44]. Datasets gathered to aid machine understanding of such hand-object interactions rely on additional instrumentation of the users’ hands [16], motion capture systems with hand-attached markers [15, 68], or multi-view camera rigs [5, 25, 38, 75, 77] to capture accurate ground-truth poses of users’ hands under the higher degree of occlusion caused by the object.

**Hand-object contact estimation** In addition to object-relative hand pose, prior work has aimed to estimate contact points between the users’ hands and external objects [15, 68]. Research has shown that when used as input proxies, real-world physical objects improve input control and provide

haptic feedback [8]. For interactive research purposes, external tracking systems [8, 58] and wearable sensors such as acoustic sensors [62] and inertial measurement units were used to estimate contact [17, 23, 50, 63, 64, 69]. Additionally, vision-based techniques have been developed that use fiducial markers [40], active illumination for shadow creation [42, 72], vibration detection [65], or depth sensing [14, 24, 73, 74]. More recent work estimates touch using passive cameras without additional instrumentation on the user’s hand or surface, enabling deployment on commercial mixed reality headsets [59, 66]. More detailed contact maps are inferred based on the intersection of tracked hand and object meshes [15, 18, 38, 56, 68, 77], requiring sub-millimeter accuracy—a challenging task for complex gestures due to soft tissue dynamics. To address this, Brahmabhatt et al. [3] used thermal imaging to obtain accurate contact maps. Additionally, prior efforts have utilized simulations to obtain more granular labels about contacting tissue [11, 29, 78].

**Hand pressure estimation** Moving beyond the mere detection of contact, prior work has estimated the pressure forces applied during hand interactions, which is crucial for robotic grasping tasks [9, 48] and provides an additional control dimension for input [57]. To estimate pressure from monocular images, visual cues such as fingernail alterations [6, 49] or surface deformations [32, 51] during press events have been used. Changes in object trajectory and interaction forces [13, 41, 55] also offer insights but are ineffective with static objects like tables and walls. Accurate pressure labels for training usually require instrumenting the user’s hands with gloves [4, 43, 67] or the surface with force sensors [19, 55, 61], ideally flexible or conforming to various shapes [2, 36, 46]. However, this alters the visual appearance and tactile features of the hands and surface, affecting interaction and limiting generalization to bare hands and uninstrumented surfaces. Grady et al. [19, 20] collected two datasets with ground-truth pressure maps using a Sensel Morph [33] pressure sensor to train a neural network for estimating contact regions on surfaces from single RGB images. However, their method relies solely on exocentric static cameras that clearly capture the fingertips.

With EgoPressure, we provide the first dataset containing egocentric and multi-view RGB-D images of a bare hand interacting with a surface, along with synchronized pressure data, hand poses, and meshes (see Table 1).

## 3. Marker-less Annotation Method

To capture accurate hand poses and meshes without markers, we developed a multi-camera hand pose annotation method using the MANO hand model [60], differentiable rendering and multi-objective optimization. Figure 2 shows an overview of our method, which relies on  $C$  static cameras

Dataset	frames	participants	hand pose	hand mesh	markerless	real	egocentric	multiview	RGB	depth	contact	pressure	
												surface	hand
EgoPressure (ours)	4.3M	21	✓	✓	✓	✓	✓	✓	✓	✓	Pressure sensor	✓	✓
ContactLabelDB [20]	2.9M	51	×	×	✓	✓	×	✓	✓	×	Pressure sensor	✓	×
PressureVisionDB [19]	3.0M	36	×	×	✓	✓	×	✓	✓	×	Pressure sensor	✓	×
ContactPose [3]	3.0M	50	✓	✓	✓	✓	×	✓	✓	✓	Thermal imprint	×	×
GRAB [68]	1.6M	10	✓	✓	×	✓	×	×	×	×	Inferred from Pose	×	×
ARCTIC [15]	2.1M	10	✓	✓	×	✓	✓	✓	✓	✓	Inferred from Pose	×	×
H2O [38]	571k	4	✓	✓	✓	✓	✓	✓	✓	✓	Inferred from Pose	×	×
OakInk [75]	230k	12	✓	✓	✓	✓	×	✓	✓	✓	Inferred from Pose	×	×
OakInk-2 [77]	4.01M	9	✓	✓	✓	✓	✓	✓	✓	×	Inferred from Pose	×	×
DexYCB [5]	582k	10	✓	✓	✓	✓	×	✓	✓	✓	Inferred from Pose	×	×
HO-3D [25]	103k	10	✓	✓	✓	✓	×	✓	✓	✓	Inferred from Pose	×	×
TACO [45]	5.2M	14	✓	✓	✓	✓	✓	✓	✓	✓	Inferred from Pose	×	×

Table 1. **Comparison between EgoPressure and selected hand-contact datasets.** The majority of prior datasets infer contacts based on hand and object pose. ContactLabelDB and PressureVisionDB also include ground-truth touch pressure but are limited to static cameras and do not provide accurate hand poses and meshes. Please see Supp. for the full table.

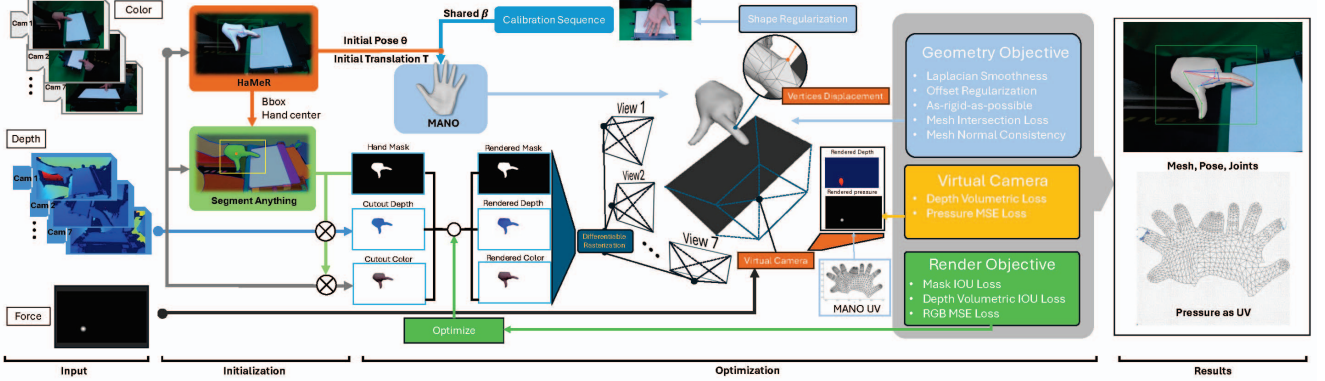


Figure 2. **Method overview.** The input for our annotation method consists of RGB-D images captured by 7 static Azure Kinect cameras and the pressure frame from a Sensel Morph touchpad. We leverage Segment-Anything [37] and HaMeR [54] to obtain initial hand poses and masks. We refine the initial hand pose and shape estimates through differentiable rasterization [7] optimization across all static camera views. Using an additional virtual orthogonal camera placed below the touchpad, we reproject the captured pressure frame onto the hand mesh by optimizing the pressure as a texture feature of the corresponding UV map, while ensuring contact between the touchpad and all contact vertices.

and a pressure-sensitive touchpad. Please see Supp. S2.2 for a detailed evaluation of our annotation method.

### 3.1. Automatic hand pose initialization

We use HaMeR [54] to estimate an initial MANO hand pose  $\theta_{\text{init}}$  and translation  $t_{\text{init}}$  for each static camera. Since HaMeR’s prediction is based on a single RGB image, there is a scale-translation ambiguity, which we resolve by triangulating the root joints from the 7 static camera views. The orientation and hand pose are then initialized based on the output of a single camera view. HaMeR also provides a bounding box, which we use along with the 2D projected hand root as input to Segment-Anything (SAM) [37], from which we obtain an annotated segmentation mask  $M_{\text{gt}}$  for the hand in each static camera image.

### 3.2. Annotation refinement

Based on the initial hand pose, we obtain refined hand pose annotations via the following optimization using the input from the  $C$  static cameras. We use the MANO [30, 60] model for mesh representation with 25 PCA components

and employ the DIB-R [7] differentiable renderer. The annotations include the hand pose  $\theta$ , hand translation  $t$ , vertex displacement  $D_{\text{vert}}$  in world coordinates, and the pressure over the hand mesh in the form of a texture map  $T_P$ . All static cameras are pre-calibrated, allowing us to project the hand mesh into the frame of each static camera  $i$  using the extrinsic parameters  $[R_{\text{cam}}^i | t_{\text{cam}}^i]$ .

**$\beta$ -calibration** For the MANO shape parameters  $\beta$ , we use separate calibration sequences for each hand of each participant, during which the participant slowly turns their hand to be visible from all cameras, with fingers spread. For these sequences, we also optimize the MANO shape parameters  $\beta$  with  $l_2$  regularization in the previous optimization. The shape parameters are then reused across all other sequences for the participant, keeping  $\beta$  fixed during subsequent optimizations.

Following HARP [35], our method consists of two stages: (1) POSE OPTIMIZATION and (2) SHAPE REFINEMENT.

In the first stage, POSE OPTIMIZATION, we annotate the hand pose  $\theta$  and translation  $t$ . The hand mesh  $\Theta$  can be derived directly from the MANO model [30], expressed



as  $\Theta = \text{MANO}(\theta, \beta) + \mathbf{t}$ . We note that certain parts of the hand, such as fingers, may not be visible from all camera angles—for instance, fingers obscured by the palm in a curled gesture. To address this, we incorporate the mesh intersection loss  $\mathcal{L}_{\text{insec}}$  [34, 70]. The objective function is then defined as:

$$\mathcal{L}_{\text{pose}}(\Theta) = \mathcal{L}_{\mathcal{R}}(\Theta) + \mathcal{L}_{\text{insec}}(\Theta). \quad (1)$$

The rendering objective  $\mathcal{L}_{\mathcal{R}}$  and the mesh intersection loss  $\mathcal{L}_{\text{insec}}$  will be detailed in Supp. S2.1.

In the SHAPE REFINEMENT stage, the pose  $\theta$  and translation  $\mathbf{t}$  of the hand remain fixed. The optimization process introduces vertex displacement  $\mathbf{D}_{\text{vert}}$ . Each vertex is adjusted by an offset along its normal vector  $\vec{n}$ , which is computed from the last epoch of the POSE OPTIMIZATION stage, to minimize the rendering loss  $\mathcal{L}_{\mathcal{R}}(\Theta^*)$ . Consequently, the refined hand mesh  $\Theta^*$  can be expressed as  $\Theta^* = \Theta + \vec{n} \cdot \mathbf{D}_{\text{vert}}$ . To ensure a reasonable mesh, the geometry objective  $\mathcal{L}_{\mathcal{G}}$  is also included in the optimization. Additionally, we introduce a virtual render  $\tilde{\mathcal{R}}$  to optimize pressure as a UV map  $\mathcal{T}_P$  and minimize the distance between the hand mesh  $\Theta^*$  and the contact area on the surface of the touchpad. The objective function  $\mathcal{L}_{\text{shape}}$  for this stage is as follows:

$$\mathcal{L}_{\text{shape}}(\Theta^*) = \mathcal{L}_{\mathcal{R}}(\Theta^*) + \mathcal{L}_{\mathcal{G}}(\Theta^*) + \mathcal{L}_{\tilde{\mathcal{R}}}(\Theta^*). \quad (2)$$

The virtual render  $\tilde{\mathcal{R}}$ , and its objective  $\mathcal{L}_{\tilde{\mathcal{R}}}$  will be explained in the next section and the other terms in the geometry objective  $\mathcal{L}_{\mathcal{G}}$  will be detailed in Supp. S2.1.2.

### 3.2.1. Virtual Render for Contact and Pressure

As shown in Figure 2, we also incorporate the captured pressure data in the optimization as a hand mesh texture feature for our proposed virtual rendering method. For this, we position a virtual orthogonal camera  $\tilde{\mathcal{R}}$  under the touchpad, oriented upwards in the world coordinate system. The render size matches the resolution of the touchpad, and the camera’s plane overlaps with the touchpad’s sensing surface. The goal is for the rendered pressure  $\tilde{\mathcal{R}}_P(\Theta^*, \mathcal{T}_P)$  on the hand mesh, with texture mapping of an optimized pressure UV map  $\mathcal{T}_P$ , to align with the input pressure  $\mathbf{P}_{\text{gt}}$ .

Additionally, we infer the contact area from  $\mathbf{P}_{\text{gt}}$  using a simple pressure threshold. Using this contact area as a mask, we ensure that the masked rendered z-axis depth  $\tilde{\mathcal{R}}_D(\Theta^*)[z]$  aligns with the distance  $Z_{v2p}$  from the camera to the touchpad, thereby ensuring physical contact.

The objective function  $\mathcal{L}_{\tilde{\mathcal{R}}}(\Theta^*)$  for the virtual render is:

$$\begin{aligned} \mathcal{L}_{\tilde{\mathcal{R}}}(\Theta^*) = & \text{MSE}(\tilde{\mathcal{R}}_P(\Theta^*, \mathcal{T}_P), \mathbf{P}_{\text{gt}}) \\ & + \left| \mathbb{I}(\mathbf{P}_{\text{gt}} > 0) \odot (\tilde{\mathcal{R}}_D(\Theta^*)[z] - Z_{v2p}) \right|_1. \end{aligned} \quad (3)$$

## 4. EgoPressure Dataset

EgoPressure comprises 4.3M RGB-D frames ( $2560 \times 1440$  for static camera,  $1920 \times 1080$  for egocentric camera) capturing interactions of both left and right hands (see Figure 7) with a touch and pressure-sensitive planar surface. The dataset features 21 participants performing 31 distinct gestures, such as touch, drag, pinch, and press, with each hand (see Figure 6). It includes a total of 5.0 hours of hand gesture footage comprised of synchronized RGB-D frames from seven calibrated static cameras and one head-mounted camera, along with ground-truth pressure maps from the pressure-sensitive surface captured at a frame rate of 30 fps. We used four different surface textures for the data capture rig, which also includes a green wall to facilitate synthetic background augmentation. Additionally, we provide high-fidelity hand pose and mesh data for the hands during interactions based on our proposed annotation method (see Section 3), as well as the tracked pose of the head-mounted camera. With EgoPressure, we offer a substantial dataset for egocentric hand pose and pressure estimation during interactions with rigid surfaces, thereby advancing machine understanding of human interaction with their surroundings through the fundamental modality of touch.

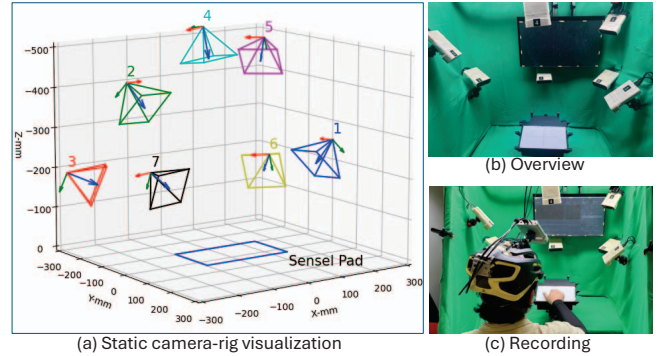


Figure 3. 7 static + 1 egocentric camera rig.

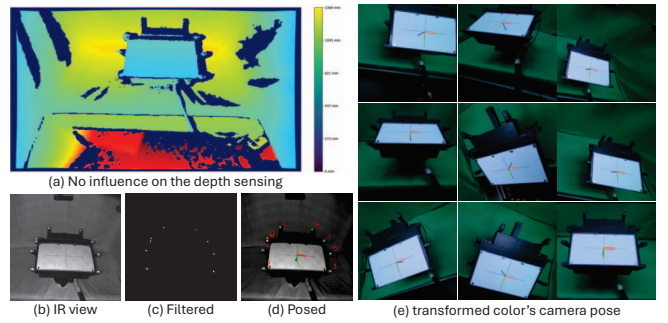


Figure 4. Camera pose tracking with IR makers.

### 4.1. Data capture setup

To capture accurate ground-truth labels for hand pose and pressure from egocentric views, we constructed a data capture rig that integrates a pressure-sensitive touchpad (Sensel Morph [33]) for touch and pressure sensing, along with

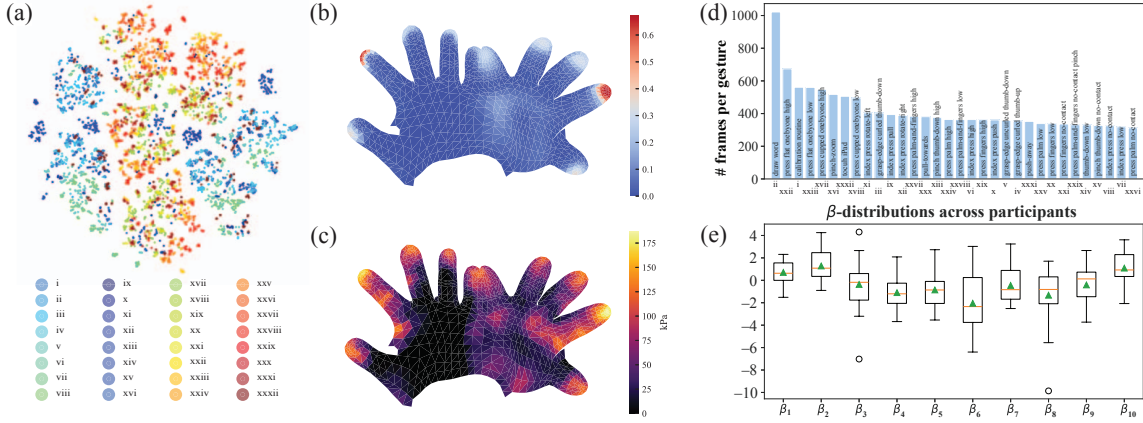


Figure 5. **Dataset Statistics** (a) t-SNE [71] visualization of hand pose frames  $\theta$  over our dataset, with color coding for different gestures. All gestures are listed in Table S5 of the Supp. (b) Ratio of touch frames with contact for each vertex. (c) Maximum pressure over hand vertices across the dataset. (d) Mean length of performed gestures. (e) Distribution of  $\beta$  values across participants.



Figure 6. **Thumbnail of different poses in egocentric views.**

seven static and one head-mounted RGB-D camera (Azure Kinect [1]) to capture RGB and depth images (see Figure 3). The touchpad (Sensel Morph), measuring  $240 \times 169.5$  mm, is mounted on a tripod head. We use four different texture overlays (white, green, dark wood, light wood) printed on paper and placed over the Sensel Morph pad across participants. The seven static Azure Kinect cameras are attached to the aluminum frame, and the head-mounted camera is fixed on a helmet. The frame also holds a computer display and is surrounded by a green screen.

All cameras and the touchpad are connected to two workstations (Intel Core i7-9700K, Nvidia GeForce RTX 3070), their timestamps are synchronized via a Raspberry Pi CM4 using PTP, which also triggers all Azure Kinect cameras simultaneously at a frame rate of 30 fps. We varied the Kinect

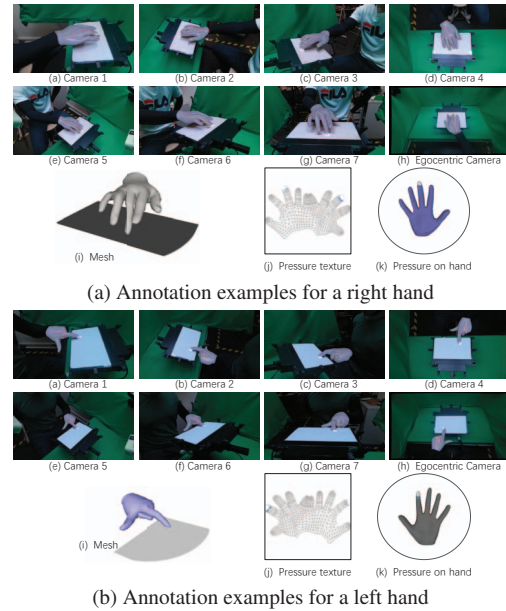


Figure 7. **Sample data from EgoPressure.**

camera exposure (2.5 ms or 10 ms) and overhead lighting in three conditions across participants: dark (2 tubes active, 2.5 ms), medium (2 tubes, 10 ms), and bright (4 tubes, 10 ms). We minimize reliance on shadows via diffuse light sources. **Head-mounted camera tracking** To obtain accurate poses of the head-mounted camera, we attach nine active infrared markers around the Sensel Morph pad in an asymmetric layout (see Figure 4). These markers, controlled by the Raspberry Pi CM4, are identifiable in the Azure Kinect’s infrared image using simple thresholding (saturating the range of values of the infrared camera). The markers are turned on simultaneously, allowing for the computation of the camera pose via Perspective-N-Points and enabling an accurate evaluation of the temporal synchronization between cameras and the touchpad (see Supp. S3.5).

## 4.2. Participants

We recruited 21 participants from our institution (6 female, 15 male, ages 23–32 years, mean age = 26 years), ensuring a broad representation to cover broad hand anatomies. Participants’ heights ranged from 160–194 cm (mean = 174, SD = 9), weights from 51–95 kg (mean = 69, SD = 14), and middle finger lengths from 7.3–9.2 cm (mean = 7.9, SD = 0.5) (see Figure 5 for distribution of MANO  $\beta$ -values).

## 4.3. Data acquisition procedure

Participants sat on an adjustable stool in front of the apparatus, wearing a helmet with a mounted camera pointing towards the Sensel Morph and a black arm sleeve on each arm up to the wrist. Before starting the data capture, the experimenter explained the task and the purpose of the study. They then signed a consent form and provided demographic information. The participants first performed a calibration gesture by slowly turning each hand, with fingers spread, within the camera rig. After calibration, participants conducted 31 different gestures, including touch, press, and drag gestures of varying strength, with each hand on the Sensel Morph touchpad (see Supp. S3.1 for a description of gestures). Each gesture was repeated 5 times if it involved a single touch action (e.g., press index finger) and 3 times if it involved a sequence of sequential touches (e.g., draw letters). Before each gesture, participants watched a video demonstrating how to perform the corresponding gesture with written instructions on a computer monitor in front of them. The experimenter guided the participants throughout the study, which took around 1 hour per participant. Participants could take a break after each gesture and received a chocolate bar as gratitude for their participation. In total, we recorded 6216 different gestures, i.e., 21 participants  $\times$  2 hands  $\times$  (1 calibration + 27  $\times$  5 + 4  $\times$  3) gestures.

## 4.4. Data statistics

The average length of each motion sequence is 14 seconds, with an almost equal balance between frames capturing the left and right hands. Figure 5 shows the mean sequence lengths across gestures. Approximately 45.1% of all frames capture the hand in contact with the pressure-sensitive pad. Figure 5b visualizes the ratio of contact frames with a given vertex touching the surface, and Figure 5c shows the maximum pressure measured for each vertex. Following Grady et al. [19], we set a threshold of 0.5 kPa as the minimum effective pressure to discard diffuse readings from the touchpad.

## 5. Benchmark Evaluation

Previous work estimates applied pressure maps using only RGB images [19, 20]. With EgoPressure, we explore the advantages of incorporating accurate hand poses as additional input, which naturally provide richer context about the inter-

action. We introduce new benchmarks for estimating hand pressure using both RGB images and 3D hand poses. Additionally, we propose a novel network architecture that jointly estimates, from a single RGB image, the pressure applied to both an external surface and across the hand, providing a deeper understanding of the regions of the hand involved throughout the interaction.

### 5.1. Image-projected Pressure Baselines

We test our hypothesis that incorporating hand pose as an additional input enhances pressure estimation. To this end, we design a straightforward extension of PressureVisionNet [19]. Specifically, we augment the original encoder-decoder segmentation architecture, which was designed for RGB inputs only, by adding an additional channel for 2.5D hand keypoints. This involves projecting the 21 3D hand joints onto the image plane and incorporating their depth (z-coordinate) from the egocentric camera’s coordinate system, scaled to millimeters.

We evaluate PressureVisionNet and the pose-augmented network on EgoPressure in three setups: (1) trained/tested on egocentric views, (2) trained/tested on the same exocentric views, and (3) trained on camera views 2,3,4,5; tested on 1,6,7. In all experiments, data from 15 participants is used for training and validation, while data from 6 participants is held out as the test set. To evaluate the augmented network, we use both the ground-truth hand joints from our annotations and the predicted hand joints from HaMeR [54]. The HaMeR-estimated hand poses serve as a fair baseline, reflecting the performance of state-of-the-art RGB-based hand pose estimators, while the ground-truth joints provide an upper bound, demonstrating the potential improvements achievable with more accurate hand poses.

The results are summarized in Table 2. Incorporating 2.5D hand joints improves performance in both egocentric and exocentric views and enhances generalization to unseen camera views. Figure 8 provides additional qualitative results, demonstrating the benefits of incorporating hand pose information. Further details on the architecture, training process, and evaluation metrics can be found in Supp. S1.1

Model	Train	Eval.	Modality	Cont. IoU $\uparrow$	Vol. IoU $\uparrow$	MAE $\downarrow$	Temp. Acc. $\uparrow$
PressureVisionNet [19]	ego.	ego.	RGB	55.73	38.64	53.60	91.68
[19] w. [54] pose	ego.	ego.	RGB & pred pose	56.25	40.52	55.23	91.67
[19] w. GT pose	ego.	ego.	RGB & GT pose	58.80	41.39	53.79	92.17
PressureVisionNet [19]	exo. (2,3,4,5)	exo. (2,3,4,5)	RGB	62.11	44.73	43.15	93.61
[19] w. [54] pose	exo. (2,3,4,5)	exo. (2,3,4,5)	RGB & pred pose	62.95	45.01	42.53	93.83
[19] w. GT pose	exo. (2,3,4,5)	exo. (2,3,4,5)	RGB & GT pose	64.39	47.58	41.72	94.18
PressureVisionNet [19]	exo. (2,3,4,5)	exo. (1,6,7)	RGB	36.82	25.05	62.22	83.40
[19] w. [54] pose	exo. (2,3,4,5)	exo. (1,6,7)	RGB & pred pose	38.46	28.10	51.50	86.34
[19] w. GT pose	exo. (2,3,4,5)	exo. (1,6,7)	RGB & GT pose	43.04	31.39	49.45	89.78

Table 2. **Image-projected pressure estimation using different inputs.** Our high-fidelity hand pose annotations improve contact IoU [%], volumetric IoU [%], MAE [Pa], and temporal accuracy [%] compared to using no hand poses or HaMeR [54] hand poses as additional input for novel exocentric and egocentric views.



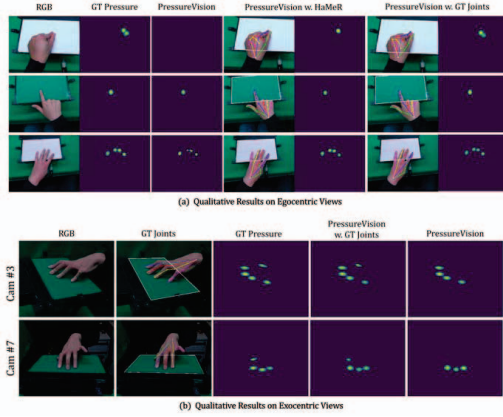


Figure 8. **Qualitative results.** We present the egocentric experiment results in Subfigure (a). In Subfigure (b), both baseline models are trained using camera views 2, 3, 4, and 5. We display the results for one seen view and one unseen view. Additionally, we overlay the 2D keypoints predicted by HaMeR [54] and our annotated ground truth on the input image. For better visualization, the contour of the touch sensing area is also highlighted as a reference.

## 5.2. First Hand-Projected Pressure Baseline

Both the original PressureVision framework [19] and its subsequent iteration, PressureVision++ [20], predict 2D hand pressure on the image plane. However, this introduces ambiguity about the exact manifestation of this pressure between hands and objects within the 3D space.

To address this, we introduce a new baseline model, *PressureFormer*, which estimates pressure as a UV map of a 3D hand mesh, enabling projection both as 3D pressure onto the hand surface and as 2D pressure onto the image space.

As illustrated in Figure 9, our model builds upon HaMeR [54]. It processes the hand vertices  $V_{hand}$  in the camera frame and the image feature tokens from HaMeR’s Vision Transformer (ViT) [12]. A transformer-based decoder receives  $V_{hand}$  as multiple input tokens while cross-attending to the image feature tokens from the ViT. Each output token represents a  $D$ -dimensional feature for a corresponding mesh vertex, which we then map onto a UV feature map using the UV coordinates of the MANO model [60]. Given the sparsity of the UV feature map post-projection, we apply two convolutional layers for neural interpolation and reduce the dimensions to the number of force classes  $C$  to predict the quantified UV-pressure map  $U_{pred}$ .

Firstly, we compute the coarse UV-pressure loss  $\mathcal{L}_c$  between  $U_{pred}$  and the ground-truth UV-pressure map  $U_{gt}$ , quantified from the scalar UV-Pressure  $\mathcal{T}_P$  of our dataset. Subsequently, we render the pressure  $P_{pred}$  back onto the original image plane using the  $M_{hand}$  mesh of vertices  $V_{hand}$  and texture mapping the predicted  $U_{pred}$  UV-pressure map. Using a differentiable renderer [7], we invert the z-normal and z-axis of the face vertices to identify the mesh faces furthest from the camera (i.e., occluded vertices) as

places of contact. This allows us to compute the pressure loss  $\mathcal{L}_p$  against the ground-truth pressure  $P_{gt}$ . Both  $\mathcal{L}_p$  and  $\mathcal{L}_c$  employ cross entropy loss. The training of PressureFormer is supervised by a loss function defined as:

$$\mathcal{L}_{PF} = w_1 \mathcal{L}_c + w_2 \mathcal{L}_p. \quad (4)$$

For comparison, we project the image-based pressure maps from PressureVisionNet and its hand-pose-augmented baseline onto the corresponding hand mesh estimated from the same image using HaMeR [54]. Similarly, this process involves identifying the hand mesh faces farthest from the camera and rasterizing the 2D pressure map onto the UV map (see Supp. Figure S2). We also evaluate a variant of PressureFormer trained without explicit UV loss supervision. We thus introduce a benchmarking task that assesses the accuracy of pressure estimation on the hand surface and the performance of jointly estimating pressure and hand mesh.

We trained our PressureFormer and baseline models using images from all camera views of 15 participants, incorporating a hand-centered crop. During training, we applied data augmentation techniques, including shifting, rescaling, and rotating. We evaluated the models on a held-out test set of six participants using (1) all camera views and (2) only egocentric camera images. Additionally, we assessed the generalization of the models to the test set of PressureVision [19]. We evaluate the accuracy of the estimated pressure map in both image space and UV space when ground truth data is available (see Supp. S1.3).

**Results** The results are summarized in Table 3. PressureFormer outperforms all image-projected pressure baselines in Contact IoU and Volumetric IoU on the UV pressure map. It also attains the highest Contact IoU on the image-projected pressure map and shows better generalization to PressureVision. The hand-pose-augmented baseline, which directly predicts pressure onto the camera image, achieves the best Volumetric IoU on the image-based pressure map. These results highlight the value of incorporating hand pose information for pressure estimation and jointly estimating hand pose and pressure for more coherent interaction modeling. Additionally, the results underline the value of the coarse UV-pressure loss in enhancing the accuracy of the pressure predictions on the UV map (see Supp. Figure S3). Figure 11 provides a qualitative comparison of the UV pressure maps estimated by the three baseline methods.

**Generalization of PressureFormer** PressureFormer employs a UV-pressure map that enhances the generalization of hand contact and pressure prediction for more complex objects. Unlike estimating pressure on the image plane, which focuses on hand-surface interactions, the UV map captures pressure on the hand vertices in 3D space. As PressureFormer uses the pretrained HaMeR [54] model as its backbone to extract hand vertices and image features from vision transformer tokens, it can effectively handle diverse

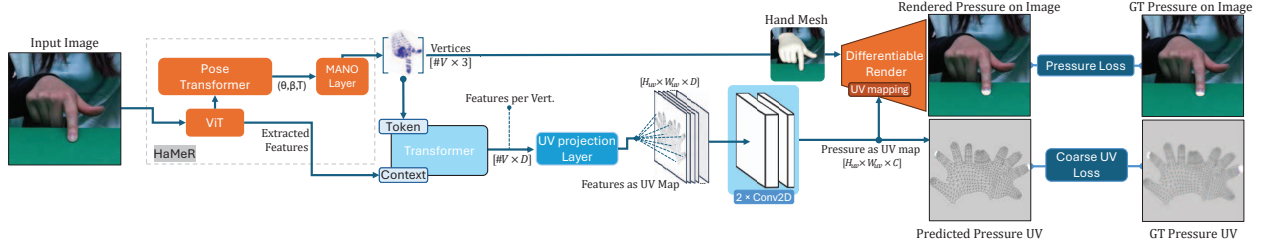


Figure 9. **PressureFormer** uses HaMeR’s hand vertices and image feature tokens to estimate the pressure distribution over the UV map. We employ a differentiable renderer [7] to project the pressure back onto the image plane by texture-mapping it onto the predicted hand mesh.

Model	Eval. Dataset	Im. Contact IoU $\uparrow$	Im. Vol. IoU $\uparrow$	Im. MAE $\downarrow$	Temp. Acc. $\uparrow$	UV Press. Contact IoU $\uparrow$	UV Press. Vol. IoU $\uparrow$
PressureVisionNet [19]	EgoPressure (ego. & exo.)	40.71	32.11	<b>44</b>	90	21.53	16.41
[19] (w/ HaMeR [54])	EgoPressure (ego. & exo.)	42.52	<b>35.40</b>	49	<b>92</b>	24.10	17.36
PressureFormer (Ours)	EgoPressure (ego. & exo.)	<b>43.04</b>	31.57	71	89	<b>33.12</b>	<b>24.54</b>
PressureFormer (w/o $\mathcal{L}_c$ )	EgoPressure (ego. & exo.)	41.27	29.57	74	88	26.24	18.61
PressureVisionNet [19]	EgoPressure (ego.)	40.65	<b>33.91</b>	<b>47</b>	<b>87</b>	26.59	19.81
PressureFormer (Ours)	EgoPressure (ego.)	<b>42.75</b>	30.57	89	83	<b>33.51</b>	<b>23.01</b>
PressureVisionNet [19]	PressureVision (exo.)	7.54	7.11	146	55	-	-
PressureFormer (Ours)	PressureVision (exo.)	<b>29.03</b>	<b>21.71</b>	<b>121</b>	<b>79</b>	-	-

Table 3. **Performance comparison** of our PressureFormer model against image-projected pressure baselines, evaluated using temporal accuracy [%], image-based pressure metrics (Image Contact IoU, Image Vol. IoU, Image MAE [kPa]), and UV map-based pressure metrics (UV Pressure IoU, UV Pressure Vol. IoU). PressureFormer demonstrates superior performance in UV pressure IoU and UV Pressure Vol. IoU, while also achieving higher scores in image-based Contact IoU. By directly predicting pressure on the UV map, PressureFormer offers advantages, enabling accurate 3D pressure reconstruction by projecting the results onto the estimated hand surface.

hand poses while integrating hand-centric image texture information. We provide qualitative results demonstrating PressureFormer’s ability to generalize to unseen camera configurations, such as the integrated passthrough sensors of the Quest 3 (see Figure 10), as well as to unseen real-world objects and environments (see Supp. Figure S5).

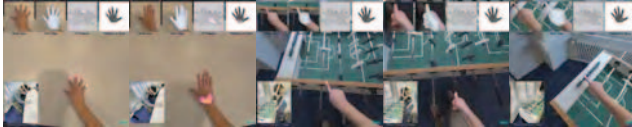


Figure 10. **PressureFormer** on data captured with Meta Quest 3.

## 6. Conclusion

**Limitations** Despite the promising generalization of the PressureFormer model, the EgoPressure dataset is limited to hand interactions with flat surfaces, as capturing precise pressure measurements on general objects without instrumenting the user’s hands remains challenging. Additionally, the dataset was collected indoors and consists only of single-hand interactions. A natural extension would be to incorporate dual-hand scenarios in more diverse environments. For a detailed discussion of these limitations, see Supp. S4).

**Summary** In this paper, we introduce EgoPressure, a novel egocentric hand pressure dataset paired with a multi-view hand pose estimation and pressure annotation method. EgoPressure includes precise 3D hand meshes, multi-view RGB and depth images, egocentric view images, and pressure intensities. We establish a new benchmark and demonstrate the effectiveness of using hand pose data in pressure estimation. Furthermore, we introduce PressureFormer, a model that directly predicts pressure on the hand mesh, along with



Figure 11. **Qualitative Results PressureFormer on our dataset.** We compare our PressureFormer with both PressureVisionNet [19] and our extended baseline model with HaMeR-estimated [54] 2.5D joint positions. Additionally, we provide visualizations of the hand mesh estimated by HaMeR, alongside the 3D pressure distribution on the hand surface derived from our predicted UV-pressure in the last two columns. Note that we transform the left-hand UV maps into the right-hand format.

relevant baselines for comparison. In conclusion, we believe EgoPressure represents an important step toward enabling machines to better understand hand-object interactions by capturing 3D pressure from an egocentric view.



## 7. Acknowledgment

We greatly appreciate all the participants who voluntarily contributed to dataset collection. We also thank Zihan Zhu, Boyang Sun, and Shaohui Liu for their insightful discussions.

## References

- [1] Azure Kinect. Azure kinect dk hardware specifications, 2019. [5](#)
- [2] Raunaq Bhirangi, Tess Hellebrekers, Carmel Majidi, and Abhinav Gupta. Reskin: versatile, replaceable, lasting tactile skins. In *5th Annual Conference on Robot Learning*, 2021. [2](#)
- [3] Samarth Brahmabhatt, Chengcheng Tang, Christopher D Twigg, Charles C Kemp, and James Hays. Contactpose: A dataset of grasps with object contact and hand pose. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16*, pages 361–378. Springer, 2020. [2](#), [3](#)
- [4] Gereon H Büscher, Risto Kõiva, Carsten Schürmann, Robert Haschke, and Helge J Ritter. Flexible and stretchable fabric-based tactile sensor. *Robotics and Autonomous Systems*, 63: 244–252, 2015. [2](#)
- [5] Yu-Wei Chao, Wei Yang, Yu Xiang, Pavlo Molchanov, Ankur Handa, Jonathan Tremblay, Yashraj S Narang, Karl Van Wyk, Umar Iqbal, Stan Birchfield, et al. Dexycb: A benchmark for capturing hand grasping of objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9044–9053, 2021. [2](#), [3](#)
- [6] Nutan Chen, Göran Westling, Benoni B Edin, and Patrick van der Smagt. Estimating fingertip forces, torques, and local curvatures from fingernail images. *Robotica*, 38(7):1242–1262, 2020. [2](#)
- [7] Wenzheng Chen, Jun Gao, Huan Ling, Edward Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. In *Advances In Neural Information Processing Systems*, 2019. [3](#), [7](#), [8](#)
- [8] Yi Fei Cheng, Tiffany Luong, Andreas Rene Fender, Paul Strelí, and Christian Holz. Comfortable user interfaces: Surfaces reduce input error, time, and exertion for tabletop and mid-air user interfaces. In *2022 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 2022. [1](#), [2](#)
- [9] Sammy Christen, Muhammed Kocabas, Emre Aksan, Jemin Hwangbo, Jie Song, and Otmar Hilliges. D-grasp: Physically plausible dynamic grasp synthesis for hand-object interactions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20577–20586, 2022. [1](#), [2](#)
- [10] Jeremy A Collins, Cody Houff, Patrick Grady, and Charles C Kemp. Visual contact pressure estimation for grippers in the wild. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10947–10954. IEEE, 2023. [1](#)
- [11] Enric Corona, Albert Pumarola, Guillem Alenya, Francesc Moreno-Noguer, and Grégory Rogez. Ganhand: Predicting human grasp affordances in multi-object scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5031–5041, 2020. [2](#)
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [7](#)
- [13] Kiana Ehsani, Shubham Tulsiani, Saurabh Gupta, Ali Farhadi, and Abhinav Gupta. Use the force, luke! learning to predict physical forces by simulating effects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 224–233, 2020. [2](#)
- [14] Neil Xu Fan and Robert Xiao. Reducing the latency of touch tracking on ad-hoc surfaces. *Proc. ACM Hum.-Comput. Interact.*, 6(ISS), 2022. [2](#)
- [15] Zicong Fan, Omid Taheri, Dimitrios Tzionas, Muhammed Kocabas, Manuel Kaufmann, Michael J Black, and Otmar Hilliges. Arctic: A dataset for dexterous bimanual hand-object manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12943–12954, 2023. [2](#), [3](#)
- [16] Guillermo Garcia-Hernando, Shanxin Yuan, Seungryul Baek, and Tae-Kyun Kim. First-person hand action benchmark with rgb-d videos and 3d hand pose annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 409–419, 2018. [2](#)
- [17] Jun Gong, Aakar Gupta, and Hrvoje Benko. Acustico: Surface tap detection and localization using wrist-based acoustic tdoa sensing. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*, page 406–419, New York, NY, USA, 2020. Association for Computing Machinery. [2](#)
- [18] Patrick Grady, Chengcheng Tang, Christopher D Twigg, Minh Vo, Samarth Brahmabhatt, and Charles C Kemp. Contactopt: Optimizing contact to improve grasps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1471–1481, 2021. [2](#)
- [19] Patrick Grady, Chengcheng Tang, Samarth Brahmabhatt, Christopher D Twigg, Chengde Wan, James Hays, and Charles C Kemp. Pressurevision: Estimating hand pressure from a single rgb image. In *European Conference on Computer Vision*, pages 328–345. Springer, 2022. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#)
- [20] Patrick Grady, Jeremy A Collins, Chengcheng Tang, Christopher D Twigg, Kunal Aneja, James Hays, and Charles C Kemp. Pressurevision++: Estimating fingertip pressure from diverse rgb images. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 8698–8708, 2024. [1](#), [2](#), [3](#), [6](#), [7](#)
- [21] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. [1](#)

- [22] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-exo4d: Understanding skilled human activity from first-and third-person perspectives. *arXiv preprint arXiv:2311.18259*, 2023. 1, 2
- [23] Yizheng Gu, Chun Yu, Zhipeng Li, Weiqi Li, Shuchang Xu, Xiaoying Wei, and Yuanchun Shi. Accurate and low-latency sensing of touch contact on any surface with finger-worn imu sensor. In *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology*, page 1059–1070, New York, NY, USA, 2019. Association for Computing Machinery. 2
- [24] Sean Gustafson, Christian Holz, and Patrick Baudisch. Imaginary phone: Learning imaginary interfaces by transferring spatial memory from a familiar device. In *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology*, page 283–292, New York, NY, USA, 2011. Association for Computing Machinery. 2
- [25] Shreyas Hampali, Mahdi Rad, Markus Oberweger, and Vincent Lepetit. Honnotate: A method for 3d annotation of hand and object poses. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3196–3206, 2020. 2, 3
- [26] Shreyas Hampali, Sayan Deb Sarkar, Mahdi Rad, and Vincent Lepetit. Keypoint Transformer: Solving Joint Identification in Challenging Hands and Object Interactions for Accurate 3D Pose Estimation. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11080–11090, Los Alamitos, CA, USA, 2022. IEEE Computer Society. 2
- [27] Shangchen Han, Beibei Liu, Randi Cabezas, Christopher D Twigg, Peizhao Zhang, Jeff Petkau, Tsz-Ho Yu, Chun-Jung Tai, Muzaffer Akbay, Zheng Wang, et al. Megatrack: monochrome egocentric articulated hand-tracking for virtual reality. *ACM Transactions on Graphics (ToG)*, 39(4):87–1, 2020. 2
- [28] Shangchen Han, Po-chen Wu, Yubo Zhang, Beibei Liu, Linguang Zhang, Zheng Wang, Weiguang Si, Peizhao Zhang, Yujun Cai, Tomas Hodan, et al. Umetrack: Unified multi-view end-to-end hand tracking for vr. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. 2
- [29] Yana Hasson, Gul Varol, Dimitrios Tzionas, Igor Kalevatykh, Michael J Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11807–11816, 2019. 2
- [30] Yana Hasson, Gül Varol, Dimitris Tzionas, Igor Kalevatykh, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *CVPR*, 2019. 3
- [31] Steven Henderson and Steven Feiner. Opportunistic tangible user interfaces for augmented reality. *IEEE Transactions on Visualization and Computer Graphics*, 16(1):4–16, 2010. 1
- [32] Wonjun Hwang and Soo-Chul Lim. Inferring interaction force from visual information without using physical force sensors. *Sensors*, 17(11):2455, 2017. 2
- [33] Sensel Inc. Sensel morph., 2024. 2, 4
- [34] Tero Karras. Maximizing parallelism in the construction of bvhs, octrees, and k-d trees. In *Proceedings of the Fourth ACM SIGGRAPH / Eurographics Conference on High-Performance Graphics*, pages 33–37. Eurographics Association, 2012. 4
- [35] Korrawe Karunratanakul, Sergey Prokudin, Otmar Hilliges, and Siyu Tang. HARP: Personalized Hand Reconstruction from a Monocular RGB Video. 2023. 3
- [36] Hong-Ki Kim, Seunggun Lee, and Kwang-Seok Yun. Capacitive tactile sensor array for touch screen application. *Sensors and Actuators A: Physical*, 165(1):2–7, 2011. 2
- [37] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023. 3
- [38] Taein Kwon, Bugra Tekin, Jan Stühmer, Federica Bogo, and Marc Pollefeys. H2o: Two hands manipulating objects for first person interaction recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10138–10148, 2021. 2, 3
- [39] Mike Lambeta, Tingfan Wu, Ali Sengul, Victoria Rose Most, Nolan Black, Kevin Sawyer, Romeo Mercado, Haozhi Qi, Alexander Sohn, Byron Taylor, et al. Digitizing touch with an artificial multimodal fingertip. *arXiv preprint arXiv:2411.02479*, 2024. 1
- [40] Minkyung Lee, Woontack Woo, et al. Arkb: 3d vision-based augmented reality keyboard. In *ICAT*, 2003. 2
- [41] Zongmian Li, Jiri Sedlar, Justin Carpentier, Ivan Laptev, Nicolas Mansard, and Josef Sivic. Estimating 3d motion and forces of person-object interactions from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8640–8649, 2019. 2
- [42] Chen Liang, Xutong Wang, Zisu Li, Chi Hsia, Mingming Fan, Chun Yu, and Yuanchun Shi. Shadowtouch: Enabling free-form touch-based hand-to-surface interaction with wrist-mounted illuminant by shadow projection. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–14, 2023. 2
- [43] PPS UK Limited. Tactileglove - hand pressure and force measurement., 2023. 2
- [44] Yunze Liu, Yun Liu, Che Jiang, Kangbo Lyu, Weikang Wan, Hao Shen, Boqiang Liang, Zhoujie Fu, He Wang, and Li Yi. Hoi4d: A 4d egocentric dataset for category-level human-object interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21013–21022, 2022. 2
- [45] Yun Liu, Haolin Yang, Xu Si, Ling Liu, Zipeng Li, Yuxiang Zhang, Yebin Liu, and Li Yi. Taco: Benchmarking generalizable bimanual tool-action-object understanding. *arXiv preprint arXiv:2401.08399*, 2024. 3
- [46] Yiyue Luo, Yunzhu Li, Pratyusha Sharma, Wan Shou, Kui Wu, Michael Foshey, Beichen Li, Tomás Palacios, Antonio Torralba, and Wojciech Matusik. Learning human–environment interactions using conformal tactile textiles. *Nature Electronics*, 4(3):193–201, 2021. 1, 2

- [47] Yiyue Luo, Chao Liu, Young Joong Lee, Joseph DelPreto, Kui Wu, Michael Foshey, Daniela Rus, Tomás Palacios, Yunzhu Li, Antonio Torralba, et al. Adaptive tactile interaction transfer via digitally embroidered smart gloves. *Nature communications*, 15(1):868, 2024. 1
- [48] Priyanka Mandikal and Kristen Grauman. Learning dexterous grasping with object-centric visual affordances. In *2021 IEEE international conference on robotics and automation (ICRA)*, pages 6169–6176. IEEE, 2021. 1, 2
- [49] Stephen A Mascaró and H Harry Asada. Measurement of finger posture and three-axis fingertip touch force using finger-nail sensors. *IEEE Transactions on Robotics and Automation*, 20(1):26–35, 2004. 2
- [50] Manuel Meier, Paul Strelí, Andreas Fender, and Christian Holz. Tapid: Rapid touch interaction in virtual reality using wearable sensing. In *2021 IEEE Virtual Reality and 3D User Interfaces (VR)*, pages 519–528. IEEE, 2021. 1, 2
- [51] Vimal Mollyn and Chris Harrison. Egotouch: On-body touch input using ar/vr headset cameras. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, pages 1–11, 2024. 2
- [52] Gyeongsik Moon, Shou-I Yu, He Wen, Takaaki Shiratori, and Kyoung Mu Lee. Interhand2. 6m: A dataset and baseline for 3d interacting hand pose estimation from a single rgb image. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*, pages 548–564. Springer, 2020. 2
- [53] Franziska Mueller, Dushyant Mehta, Oleksandr Sotnychenko, Srinath Sridhar, Dan Casas, and Christian Theobalt. Real-time hand tracking under occlusion from an egocentric rgb-d sensor. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1154–1163, 2017. 2
- [54] Georgios Pavlakos, Dandan Shan, Ilija Radosavovic, Angjoo Kanazawa, David Fouhey, and Jitendra Malik. Reconstructing hands in 3D with transformers. In *CVPR*, 2024. 2, 3, 6, 7, 8
- [55] Tu-Hoa Pham, Nikolaos Kyriazis, Antonis A Argyros, and Abderrahmane Kheddar. Hand-object contact force estimation from markerless visual tracking. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2883–2896, 2017. 2
- [56] Chandradeep Pokhariya, Ishaan Nikhil Shah, Angela Xing, Zekun Li, Kefan Chen, Avinash Sharma, and Srinath Sridhar. Manus: Markerless grasp capture using articulated 3d gaussians. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2197–2208, 2024. 2
- [57] Philip Quinn, Wenxin Feng, and Shumin Zhai. Deep touch: Sensing press gestures from touch image sequences. *Artificial Intelligence for Human Computer Interaction: A Modern Approach*, pages 169–192, 2021. 2
- [58] Mark Richardson, Matt Durasoff, and Robert Wang. Decoding surface touch typing from hand-tracking. In *Proceedings of the 33rd annual ACM symposium on user interface software and technology*, pages 686–696, 2020. 2
- [59] Mark Richardson, Fadi Botros, Yangyang Shi, Pinhao Guo, Bradford J Snow, Linguang Zhang, Jingming Dong, Keith Vertanen, Shugao Ma, and Robert Wang. Stegotype: Surface typing from egocentric cameras. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, pages 1–14, 2024. 2
- [60] Javier Romero, Dimitrios Tzionas, and Michael J. Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 2017. 2, 3, 7
- [61] Pressure Mapping Sensors. Tekscan., 2024. 2
- [62] Yilei Shi, Haimo Zhang, Jiashuo Cao, and Suranga Nanayakkara. Versatouch: A versatile plug-and-play system that enables touch interactions on everyday passive surfaces. In *Proceedings of the Augmented Humans International Conference*, New York, NY, USA, 2020. Association for Computing Machinery. 2
- [63] Yilei Shi, Haimo Zhang, Kaixing Zhao, Jiashuo Cao, Mengmeng Sun, and Suranga Nanayakkara. Ready, steady, touch! sensing physical contact with a finger-mounted imu. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 4(2), 2020. 2
- [64] Paul Strelí, Jiaxi Jiang, Andreas Fender, Manuel Meier, Hugo Romat, and Christian Holz. Tapttype: Ten-finger text entry on everyday surfaces via bayesian inference. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, New York, NY, USA, 2022. Association for Computing Machinery. 2
- [65] Paul Strelí, Jiaxi Jiang, Juliette Rossie, and Christian Holz. Structured light speckle: Joint ego-centric depth estimation and low-latency contact detection via remote vibrometry. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–12, 2023. 2
- [66] Paul Strelí, Mark Richardson, Fadi Botros, Shugao Ma, Robert Wang, and Christian Holz. Touchinsight: Uncertainty-aware rapid touch and text input for mixed reality from ego-centric vision. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, pages 1–16, 2024. 1, 2
- [67] Subramanian Sundaram, Petr Kellnhofer, Yunzhu Li, Jun-Yan Zhu, Antonio Torralba, and Wojciech Matusik. Learning the signatures of the human grasp using a scalable tactile glove. *Nature*, 569(7758):698–702, 2019. 2
- [68] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. Grab: A dataset of whole-body human grasping of objects. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 581–600. Springer, 2020. 2, 3
- [69] Ryo Takahashi, Masaaki Fukumoto, Changyo Han, Takuya Sasatani, Yoshiaki Narusue, and Yoshihiro Kawahara. Telemetry: A batteryless and wireless ring-shaped keyboard using passive inductive telemetry. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*, page 1161–1168, New York, NY, USA, 2020. Association for Computing Machinery. 2
- [70] Dimitrios Tzionas, Luca Ballan, Abhilash Srikantha, Pablo Aponte, Marc Pollefeys, and Juergen Gall. Capturing hands in action using discriminative salient points and physics simulation. *International Journal of Computer Vision (IJCV)*, 118(2):172–193, 2016. 4
- [71] Laurens Van der Maaten and Geoffrey Hinton. Visualizing



- data using t-sne. *Journal of machine learning research*, 9(11), 2008. [5](#)
- [72] Andrew D Wilson. Playanywhere: a compact interactive tabletop projection-vision system. In *Proceedings of the 18th annual ACM symposium on User interface software and technology*, pages 83–92, 2005. [1](#), [2](#)
  - [73] Andrew D. Wilson. Using a depth camera as a touch sensor. In *ACM International Conference on Interactive Tabletops and Surfaces*, page 69–72, New York, NY, USA, 2010. Association for Computing Machinery. [2](#)
  - [74] Robert Xiao, Julia Schwarz, Nick Throm, Andrew D. Wilson, and Hrvoje Benko. Mrtouch: Adding touch input to head-mounted mixed reality. *IEEE Transactions on Visualization and Computer Graphics*, 24(4):1653–1660, 2018. [2](#)
  - [75] Lixin Yang, Kailin Li, Xinyu Zhan, Fei Wu, Anran Xu, Liu Liu, and Cewu Lu. Oakink: A large-scale knowledge repository for understanding hand-object interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20953–20962, 2022. [2](#), [3](#)
  - [76] Shanxin Yuan, Qi Ye, Bjorn Stenger, Siddhant Jain, and Tae-Kyun Kim. Bighand2. 2m benchmark: Hand pose dataset and state of the art analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4866–4874, 2017. [2](#)
  - [77] Xinyu Zhan, Lixin Yang, Yifei Zhao, Kangrui Mao, Hanlin Xu, Zenan Lin, Kailin Li, and Cewu Lu. Oakink2: A dataset of bimanual hands-object manipulation in complex task completion. *arXiv preprint arXiv:2403.19417*, 2024. [2](#), [3](#)
  - [78] Zehao Zhu, Jiashun Wang, Yuzhe Qin, Deqing Sun, Varun Jampani, and Xiaolong Wang. Contactart: Learning 3d interaction priors for category-level articulated object and hand poses estimation. *arXiv preprint arXiv:2305.01618*, 2023. [2](#)
  - [79] Christian Zimmermann, Duygu Ceylan, Jimei Yang, Bryan Russell, Max Argus, and Thomas Brox. Freihand: A dataset for markerless capture of hand pose and shape from single rgb images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 813–822, 2019. [2](#)
  - [80] Lara Zlokapa, Yiyue Luo, Jie Xu, Michael Foshey, Kui Wu, Pulkit Agrawal, and Wojciech Matusik. An integrated design pipeline for tactile sensing robotic manipulators. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 3136–3142. IEEE, 2022. [1](#)