
Towards A Richer 2D Understanding of Hands at Scale

Anonymous Author(s)

Affiliation

Address

email

1 **Suggested reading order for reading the supplement.** The supplement is large so that it can cover
2 all information that is needed or was promised. We have ordered our supplement in a suggested
3 reading order that will be of the most interest for the casual reader. This is not to diminish the
4 importance of the other sections (e.g., the datasheet or full instructions), but these are references
5 rather than something that can be easily digested in one sitting.

6 **Correction to Numbers.** The *quantitative* results for a nearly identical model were inadvertently
7 reported in two tables in the main paper. In particular, the reported numbers come from an identical
8 model to the one described in the paper whose segments came from an internal SAM-like system
9 (Section D) as opposed to SAM [10]. The model that is shown throughout in qualitative results and
10 in other Tables is trained with SAM [10] masks. The two differ only in where their segmentation
11 ground-truth come from: the architecture, code, and all other aspects are identical.

12 The authors regret this error and do not believe that the error alters the conclusions of the paper
13 (indeed the models have nearly identical bounding box performance), but wish to report the correct
14 numbers. The two corrections are as follows.

15 *Segmentation (Table 3).* The true quantitative performance for segmentation is $\approx 15\%$ higher than
16 were reported in Table 3 of the main paper. Our submitted numbers were taken from the model trained
17 on our own SAM-like outputs. The segmentation performance should read as follows: Hand: 73.3
18 (not 55.0); Object: 51.0 (not 36.7); and Second Object: 31.9 (not 21.3). Bounding box performance
19 is effectively identical.

20 *Blur-vs-unblur Experiment (Table 5).* We reported bounding box performance here, using models
21 trained with masks from our own internal SAM-like system. as shown in Section D, bounding box
22 performance is effectively identical (Hand AP: 83.5 vs 83.5) or within the margin of error for a
23 random seed (object: 59.5 vs 58.6; second object: 44.4 vs 45.2).

24 Contents

25	A Additional Qualitative Examples	4
26	A.1 Failure Cases	4
27	A.2 Grasp Type Ranking on 4 Subsets and Ego4D	4
28	A.3 More Hand Configurations	5
29	A.4 More Qualitative Examples on Hands23	6
30	A.5 More Qualitative Examples on Ego4D	6
31	B Additional Quantitative Experiments	9
32	B.1 Extended Table of datasets.	9
33	B.2 COCO evaluation numbers for detection	9

34	B.3 Full Blur No Blur Tables	9
35	B.4 Audit for Differences in Performance across Skin Tone and Gender Presentation	9
36	C Model Architecture and Training Details	11
37	C.1 Model Architecture	11
38	C.2 Training Details	12
39	C.3 Finetuning on Ego4D	12
40	D Masks from an Internal SAM-like System	14
41	D.1 Model Architecture and Training Details	14
42	D.2 Computational Requirements	14
43	D.3 Performance and Discussion	14
44	E Data Processing and Redaction	16
45	E.1 Face Blurring	16
46	E.2 Child Detection	17
47	F New Videos	17
48	F.1 Video Selection	17
49	F.2 Frame Selection	18
50	F.3 Search Grammar	19
51	G Datasheet for Hands23	21
52	H Data Annotation and Instructions	30
53	H.1 Hand Detection	31
54	H.1.1 Annotation for VISOR	31
55	H.1.2 Annotation for Articulation, New Videos	31
56	H.1.3 Annotation for COCO	32
57	H.2 Hand Contact State	33
58	H.2.1 VISOR, Articulation, New Videos	33
59	H.2.2 COCO	33
60	H.3 Additional Annotations – Checking Hands Labeled by One Annotator	34
61	H.4 Object Box	36
62	H.5 Object Tool/Container Status	37
63	H.6 Second Box	39
64	H.7 Grasp	40
65	H.8 Prehensile-vs-Non-Prehensile Grasps	40
66	H.9 Differentiating Prehensile Grasps	40
67	H.10 Video Identification	43
68	H.11 Filtering Videos	43
69	H.12 Counting Hands	44

70	H.13 Identifying Frame Types	44
71	H.14 Image Redaction	46
72	H.14.1 Instructions for Unblurred Face Spotting	46
73	H.14.2 Instructions for Unblurred Face Bounding	47
74	H.14.3 Instructions for Spotting Minors	47
75	H.15 Polygon Labeling	49
76	H.15.1 Hands	49
77	H.15.2 Objects	50

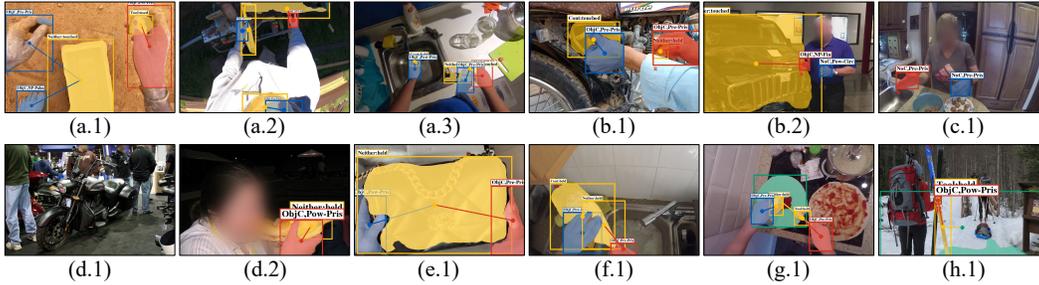


Figure 1: Failure cases. Common failure modes include false positive detection, false negative detection, missing detection, confusion between grasp predictions, confusion between fine-grained contact predictions, occlusions, and mask prediction for large scene objects.

78 A Additional Qualitative Examples

79 A.1 Failure Cases

80 Our model is not perfect. Here we discuss common failure cases that happen during our experiments
 81 as thoroughly as possible in Fig 1 ranging from common detection difficulties to hard scenarios
 82 (occlusions, shadows, super large to tiny objects, etc.).

- 83 (a) Hand false positive detection. For example, (a.1) (a.2), and (a.3) show that foot, face, or
 84 shadow are predicted as hands.
- 85 (b) Object false positive detection. In (b.1), when the hand is curled, the model hallucinates that
 86 there is an object for the right hand of the person. In (b.2), the right hand of the person is
 87 wrongly predicted as being in contact with the car.
- 88 (c) Object false negative detection. As in (c.1) the object in the left hand is in motion and has
 89 blur and is missed by the detector.
- 90 (d) Missing hand detection. In (d.1) and (d.2), the hands are small or occluded which leads to
 91 missing detection.
- 92 (e) Confusion between Pre-Pris (Precision Prismatic) and Pow-Pris (Power Prismatic). In (e.1),
 93 the two hands are holding the object with similar grasp but are predicted differently as
 94 Pre-Pris and Pow-Pris.
- 95 (f) Confusion between Pre-Pris and NP-Finger (Non-Prehensile Fingers Only). In (f.1) the right
 96 hand is predicted as Pre-Pris but it is only contacting the bottle with fingers.
- 97 (g) Confusion between tool-held and tool-used. In (g.1) the spoon is in contact with the pizza,
 98 but is predicted as being held, as opposed to used.
- 99 (h) Bad mask prediction for large scene objects. In (h.1), the ski pole (first object) is sticking
 100 into the snow (second object). Mask prediction on such kind of large scene objects is very
 101 hard.

102 A.2 Grasp Type Ranking on 4 Subsets and Ego4D

103 Understanding hand grasp is about understanding how hands and objects contact each other during
 104 hand-object interaction. It is critical for understanding the inter-relationship between hands and
 105 objects as well as transferring hand grasp manipulation ability to robot grasp manipulation. There are
 106 a lot of video data capturing tons and tons of hand activities, but how could we get useful information
 107 from the data?

108 One application of our model is to search for certain grasps in the data. Here we run our model on
 109 Hands23 testset and Ego4D valset to get all the hand grasp predictions. Ranking the grasp scores
 110 on each grasp type gives us the most typical hand grasp of each type. In Fig 2, we show the most
 111 confident sample of each grasp type on each subset of Ego4D. Since there is very little training data
 112 for Lateral grasps, no hand has lateral predicted as the most likely grasp and we therefore do not

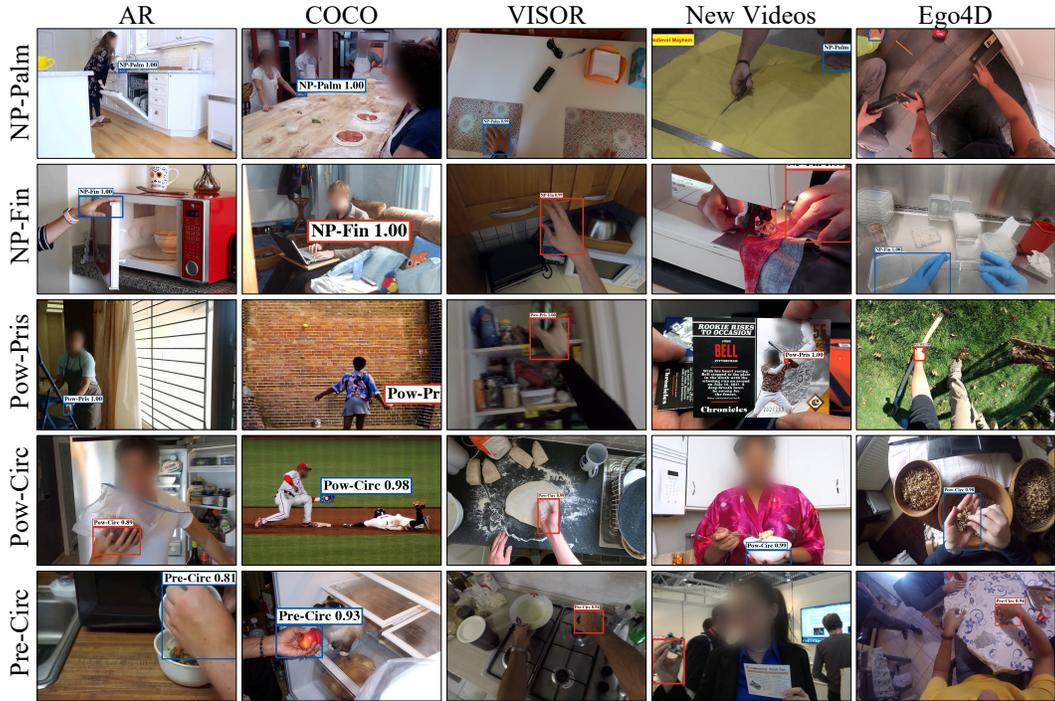


Figure 2: Grasp type ranking. Showing the Top 1 ranked sample of each grasp type on the 4 subsets (AR: Internet Articulation data [16]; COCO: COCO 2017 train [11]; VISOR [6]; and New Videos) plus Ego4D.

113 include it. We believe that our model is able to serve as a useful tool to collect hand grasp information
 114 on wild data, such as retrieving certain grasp types.

115 A.3 More Hand Configurations

116 Another application is to find various hand configurations. Previously, most hand-object interaction
 117 research is focusing on one-hand-one-object interaction. However, in reality, there are more chal-
 118 lenging hand object configurations such as bi-manual manipulation (one object interacting with two
 119 hands) and hand-tool-object interaction (the hand interacting with a tool and the tool affects the end
 120 object).

121 We present 6 interesting hand configurations here (although also note that there are more to explore),
 122 including the one $Hand \rightarrow Object \rightarrow 2nd\ Object \leftarrow Hand$ mentioned in the paper. First, we give a
 123 description of them.

- 124 • $Hand \rightarrow Object \leftarrow Hand$ (HOH) is two hands interacting with the same object.
- 125 • $Hand \rightarrow Object, Object \leftarrow Hand$ (HO, OH) is two hands interacting with different objects.
- 126 • $Hand \rightarrow Object \rightarrow 2nd\ Object \leftarrow Hand$ (HTOH) is one hand interacting with an object
 127 (tool-1) which also is interacting with a second object, while the other hand is interacting
 128 with the second object too.
- 129 • $Hand \rightarrow Object \rightarrow 2nd\ Object, Object \leftarrow Hand$ (HTO, OH) is one hand interacting with
 130 an object (tool-1) which also interacting with a second object, while the other hand is
 131 interacting with another object.
- 132 • $Hand \rightarrow Object \rightarrow 2nd\ Object \leftarrow Object \leftarrow Hand$ (HTOTH) is one hand interacting with
 133 an object (tool-1) which is interacting with a second object (2nd-obj-1), while the other
 134 hand is interacting with an object (tool-2) which is interacting with the same second object
 135 (2nd-obj-1).

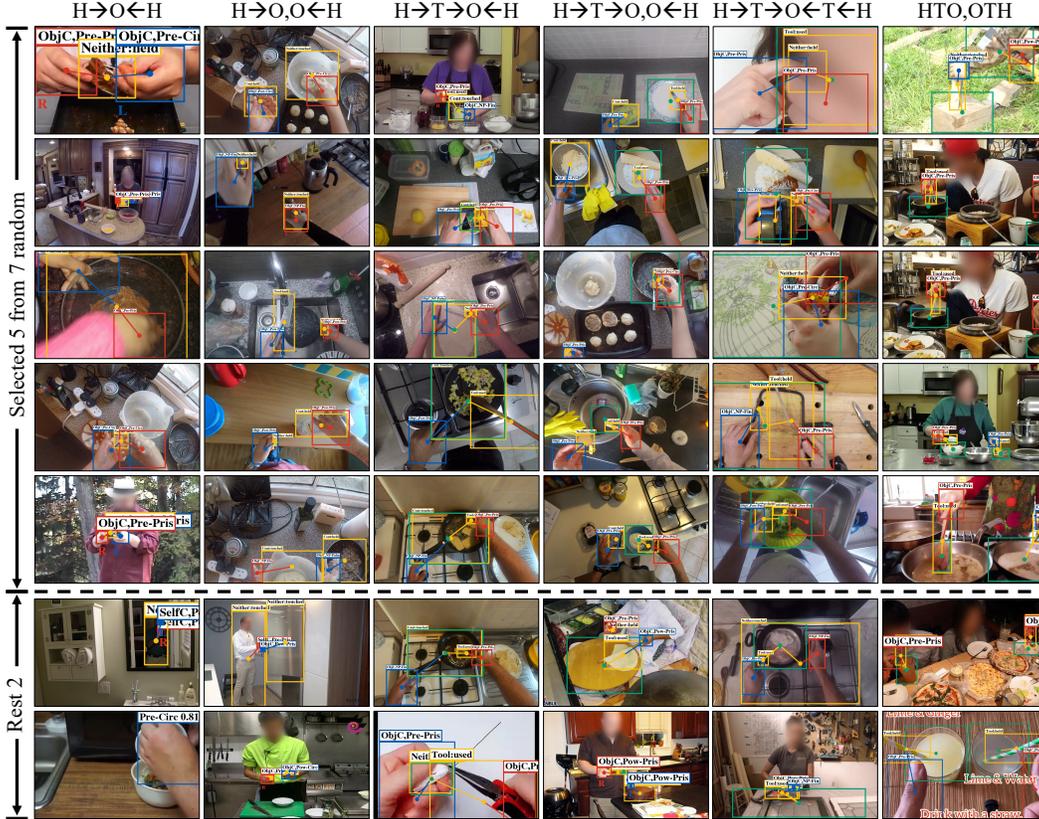


Figure 3: More hand configurations. There are various hand configurations when using hands. Here we show 6 hand configurations of bi-manual manipulations. Our model enables finding the hand-tool-object configuration.

- 136 • *Hand* → *Object* → *2nd Object*, *2nd Object* ← *Object* ← *Hand* (HTO, OTH) is one hand
 137 interacting with an object (tool-1) which is interacting with a second object (2nd-obj-1),
 138 while the other hand is interacting with an object (tool-2) which is interacting with a different
 139 second object (2nd-obj-2).

140 In Fig 3, we assume that the two hands belong to the same person. But as shown the in random 7
 141 results, there are examples (e.g. the one at row 2 col 6) of the two hands belonging to two people. In
 142 the future, incorporating [14] in the searching algorithm will help associate hands with bodies and
 143 thus make sure the two hands belong to the same body. When deciding if the two objects are the
 144 same, we threshold the bounding box IoU with a relatively high value of 0.8.

145 A.4 More Qualitative Examples on Hands23

146 In Fig 4, we provide more visualization of predictions on random images from Hands23 testset.

147 A.5 More Qualitative Examples on Ego4D

148 In Fig 5, we provide more visualization of predictions on random images from Ego4D valset.

Table 1: Extended dataset comparison. Compared with existing datasets, Hands23 is the first dataset that includes the second object annotations.

	Source	#Img	>2 Hand	#Hands	w/Obj	w/2nd	#Obj	#2nd
100DOH [18]	Video	100K	9.5%	190K	74.1%	0.0%	140K	0K
Hands23	Both	257K	8.6%	401K	71.7%	4.9%	288K	19K
<i>New Videos</i>	Video	96K	4.4%	121K	77.0%	8.2%	93K	10K
VISOR [6]	Video	38K	0.0%	58K	83.8%	10.7%	49K	6K
COCO [11]	Image	45K	33.4%	123K	64.4%	1.9%	79K	2K
<i>Artic.</i> [16]	Video	76K	3.6%	97K	67.0%	0.9%	65K	0K
VLOG [7]	Video	5K	6.5%	26.1K	-	-	-	-
VIVA [1]	Capture	5.5K	30.5%	13.2K	-	-	-	-
Ego [2]	Capture	4.8K	73.8%	15K	-	-	-	-
VGG [13]	Flickr, TV	2.7K	28.4%	4.2K	-	-	-	-
TV-Hand [15]	TV	9.5K	10.7%	8.6K	-	-	-	-
COCO-Hand [15]	Flickr	26.5K	18.4%	45.7K	-	-	-	-

Table 2: The proposed model’s performance on detection and segmentation, comparing AP50 (commonly used in past hand detection settings) and COCO AP.

	Detection (AP)			Segmentation (AP)		
	Hand	Object	2nd Object	Hand	Object	2nd Object
AP50	83.5	59.5	44.4	73.3	51.0	31.9
COCO	58.4	35.1	28.5	54.8	32.4	15.5

149 B Additional Quantitative Experiments

150 B.1 Extended Table of datasets.

151 We compare Hands23 with more existing datasets in Table 1. Hands23 is the first dataset that
 152 introduces second object annotation for understanding hand-object interaction. 100DOH also has
 153 first object annotation but the amount of first object box in Hands23 is around twice as much as that
 154 of 100DOH. VLOG has object annotation but that is clip-wise object category label instead of object
 155 bounding box label.

156 B.2 COCO evaluation numbers for detection

157 We also report the COC mAP (averaged over IOU thresholds) of the detection performance of our
 158 model in Table 2. Performance is lower, suggesting that there is lots of room for improvement by
 159 subsequent models in precise segmentation of the objects.

160 B.3 Full Blur No Blur Tables

161 We report the full performance for all four combinations of training/testing on blurred and non-blurred
 162 images in Table 3. Performance is largely identical. Training on unblurred data and testing on blurred
 163 data produces the worst results consistently; however, the gap is relatively small. We do observe a
 164 small number of false positives on faces when training on blurred data and testing on unblurred data.

165 B.4 Audit for Differences in Performance across Skin Tone and Gender Presentation

166 We report performance across Female/Male and darker skin (Fitzpatrick 4 - 6) and lighter skin
 167 (Fitzpatrick 1-3). We quantify this with both the rate (i.e., number of false detections per image as a
 168 percent) and then Fisher’s exact test. We test whether there is a difference in the number of images
 169 with an error per column. We obtained these results by selecting 100 images for each category, then
 170 excluding ambiguous cases and EPIC-KITCHENS due to its substantial skew in skin-tone. Two

Table 3: Complete comparison of training/testing on ✓blurred and ✗non-blurred scenes. Performance is largely the same across the conditions.

	Detection (AP)			Segmentation (AP)			State (Acc)			
	Hand	Obj.	2nd Obj.	Hand	Obj.	2nd Obj.	Side	Cont.	Fine	Grasp
✓Blur → ✓Blur	83.5	58.6	45.2	55.0	36.7	21.3	95.7	83.7	63.9	54.1
Blur → ✗Not Blur	83.4	58.5	45.0	54.4	36.1	21.2	95.4	83.6	63.6	54.1
✗Not Blur → Blur	82.5	57.9	43.9	54.2	35.9	20.6	94.8	83.7	63.0	53.4
✗Not Blur → ✗Not Blur	84.4	59.2	44.3	54.3	36.1	20.9	95.6	83.8	63.7	53.7

171 authors independently evaluated the outputs and counted false positives/negatives; if any annotator
 172 spotted an error, it was counted as an error.

173 There is not the yawning gap exhibited in the GenderShades [3] paper and error rates are typically
 174 quite close. There is a slight increase in errors that is not statistically significant, especially for object
 175 FPRs. We believe that various uncontrolled statistical biases are still present in the data, for instance
 176 in terms of the subject matter and composition. However, we urge that downstream users monitor
 177 output to see if they see these performance differences play out in their own data.

Table 4: Audit of performance across skin tone and presented gender. We report (in percentage) the false positive rate and false negative rate for hands and objects. We additionally report the p-value of Fisher’s exact test, testing whether there is a difference in the number of images with a false positive/negative.

	Hand FPR	Hand FNR	Obj FPR	Obj FNR
Female	4.3 ± 2.1	12.8 ± 3.8	24.5 ± 4.9	12.8 ± 3.8
Male	2.2 ± 1.6	12.4 ± 3.5	16.9 ± 4.3	10.1 ± 3.2
Fisher’s Exact p	0.68	1.00	0.27	0.81
Fitzpatrick 1-3	3.0 ± 2.1	13.6 ± 4.2	12.1 ± 4.0	9.1 ± 3.5
Fitzpatrick 4-6	3.2 ± 1.8	12.8 ± 3.8	16.0 ± 3.8	11.7 ± 3.9
Fisher’s Exact p	1.00	0.81	0.65	1.00

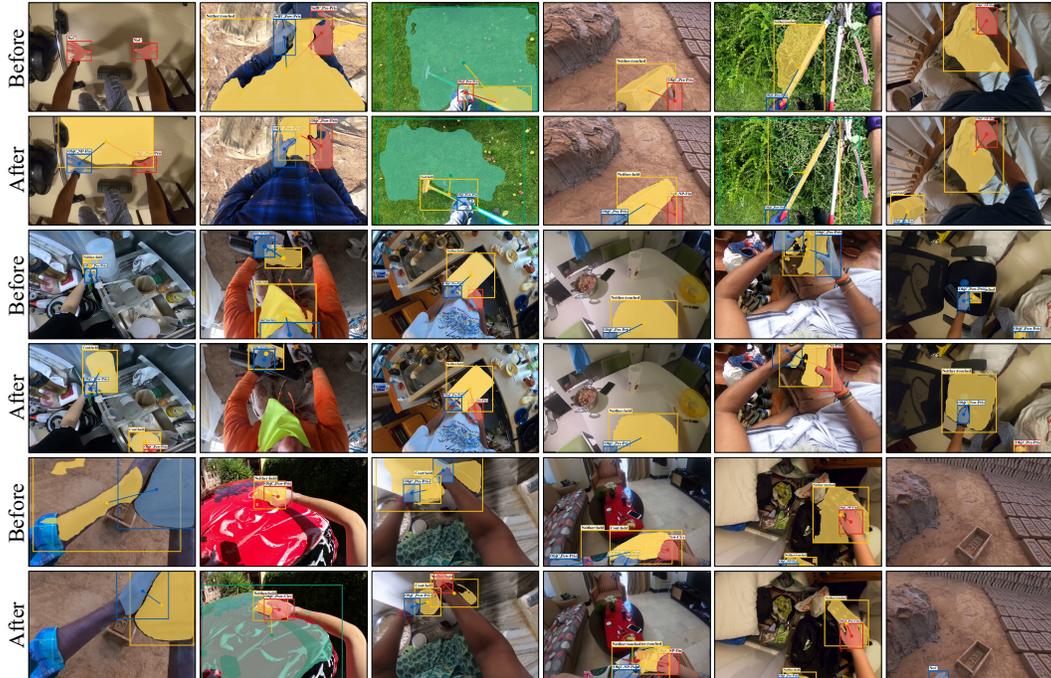


Figure 7: Performance Comparison for the model before and after finetuning on Ego4D.

Table 5: Finetuning on Ego4D. The performance on hand box, hand segmentation and hand side prediction improve after finetuning on Ego4D data. Note that Ego4D only provides side information.

	Detection (AP) Hand	Segmentation (AP) Hand	State (Acc) side
Before	86.9	60.0	0.92
After	90.9	66.3	0.97

218 C.2 Training Details

219 We have trained three different models on different versions of the data. In addition to SAM masks,
 220 we also trained our own mask segmentation models to get the automatically generated masks, which
 221 we describe in Section D. For data privacy purposes, we blurred all faces in the images. We trained
 222 models on both blurred and raw images to assess the impact of blurring.

- 223 • Model 1 (our final setting): trained on blurred images with SAM masks plus all other labels.
- 224 • Model 2: trained on blurred images with self-trained masks plus all other labels.
- 225 • Model 3: trained on raw images with self-trained masks plus all other labels.

226 The training recipe is the same for all models. The base learning rate is 0.01, which was scaled by
 227 0.1 at iterations 210000 and 250000. The models were trained for 400K iterations using 8 NVIDIA
 228 GeForce GTX 1080 Ti GPUs for around 1 week.

229 The detection loss and segmentation loss remain the same as in Detectron2. The losses for all auxiliary
 230 heads and the interaction association are scaled by 0.1; we apologize for the incorrect scaling reported
 231 in the main paper. From our previous training experiments, training the model using learning rate of
 232 0.01 without loss scaling on all auxiliary heads is unstable and leads to training divergence.

233 C.3 Finetuning on Ego4D

234 We finetuned Model 2 on Ego4D training set for another 110K iterations. Comparison results are in
 235 Fig 7 and Table 5

236 Ego4D only provides hand boxes and handside labels. In addition, we provide pseudo-labels (hand
237 contact, fine-grained object contact and hand grasp) generated automatically from our model plus
238 SAM masks for the finetuning.

239 The performance before finetuning on Ego4D shows the strong generalization ability of our model.
240 After finetuning, the performance improved. This shows that the model's performance gains with
241 finetuning on unseen data.

Table 6: The proposed model’s performance on detection and segmentation (AP), comparing using SAM [10] and our own internal system. Bounding box performance is largely identical, but segmentation is better using SAM.

	Detection (AP)			Segmentation (AP)		
	Hand	Object	2nd Object	Hand	Obj.	2nd Object
Trained on SAM [10]	83.5	59.5	44.4	73.3	51.0	31.9
Internal Masks	83.5	58.6	45.2	55.0	36.7	21.3

242 D Masks from an Internal SAM-like System

243 The masks produced by our system are automatically generate. In the paper, we use masks that come
 244 from SAM [10]. However, during the development of the project, we had developed an in-house
 245 SAM-like system. After the release of SAM, we switched to SAM. However, we document this
 246 model and its performance as an illustration of an alternate approach, since it was inadvertently used
 247 in a few tables, and to accurately capture the compute time used in the project.

248 D.1 Model Architecture and Training Details

249 Our model used aimed to automatically generate masks from available bounding boxes and used
 250 an HRNet [19] network with ResNet50 [9] backbone pretrained on ImageNet [17]. This model is
 251 trained on supervised examples that come from VISOR [6] and COCO [11].

252 To train a segmentation model in which the objective function is focused on maximizing the seg-
 253 mentation performance rather than localization ability, we crop and fixate hands and objects to the
 254 image center and pad to VISOR’s image size to have a constant resolution for training. The same
 255 preprocessing is applied to ground truth masks so that there is a pixel-to-pixel correspondence.

256 **Hands.** For hands, we crop images and masks along annotated VISOR [6] bounding boxes, pad
 257 the crops to VISOR’s image size and use available VISOR masks as our ground truth. After the
 258 first round of training on VISOR masks only, we use this model to generate pseudolabels for all
 259 other subsets. We then select good quality pseudolabels that cover greater than 70 percent of the
 260 bounding box area to be added to the training set in subsequent rounds. We repeat this process till
 261 the performance trajectory levels off. The final trained model is used to automatically generate hand
 262 masks for all subsets excluding VISOR.

263 **Objects.** Similarly for objects, we train an object segmentation model using available masks from
 264 VISOR [6] and COCO [11]. We pad all crops to VISOR’s image size and objects greater or smaller
 265 than this size are either scaled up or down during inference. In this case, we do not see an increase in
 266 performance when training on additional pseudolabels, hence we halt training after the first round.
 267 This trained model is finally used to generate masks for all images with no corresponding ground
 268 truths.

269 D.2 Computational Requirements

270 Both hand and object segmentation models were trained using a NVIDIA GeForce RTX 2080 Ti with
 271 a batch size of 1. The models were trained for a single epoch which takes about four to six days. We
 272 estimate that during development, we trained on the order of 25 versions of the model.

273 D.3 Performance and Discussion

274 Despite the distributional shift encountered when generating masks using models trained on mostly
 275 egocentric data, we found that this approach performed quite well, with some caveats.

276 **Performance.** We report performance in Table 6, using SAM outputs as ground-truth. Bounding box
 277 performance is largely unaffected by changing the labels. SAM produces better segments, by about
 278 $\approx 15\%$. For metrics and a discussion of the use of SAM outputs as ground-truth, please see the main
 279 paper’s metrics section.

280 Early in the project, when testing a model trained on only egocentric hand data, we observed that
281 the model uses a shortcut of learning to identify skin surfaces. Since egocentric data only includes
282 the camera wearer’s hands and arms, there are no negative examples of skin surfaces the model can
283 learn from. Based on this observation, we trained the on masks bounded by annotated hand boxes.
284 However, this implies that the model is highly sensitive to skin tone and can only identify hands of
285 the demographics it was trained on.

286 The object segmentation model performs relatively well on images where object boundaries are
287 clearly defined and foreground-background contrast is high. We notice a dip in performance when
288 it comes to large objects due to the limited number of large objects in the training set. Adjacently,
289 having COCO masks in our training set significantly improved performance on very small objects.

290 E Data Processing and Redaction

291 We have made substantial efforts to blur all faces and remove all children from our dataset. We now
292 describe how we did each step. All specifics about the annotation instructions appear in Section H.

293 E.1 Face Blurring

294 We followed a process that aims to blur all the recognizable faces in the dataset. This follows a
295 multi-step process that is partially automated but has several manual checks.

296 Generating boxes and masks

297 **Step 1 – Initial Boxes.** Our first round boxes come from the AWS Rekognition service. This finds
298 most of the faces in the dataset but is imperfect (hence our multiple manual steps).

299 **Step 2 – Verification.** We apply our face blurring algorithm below and then ask workers to check
300 that all faces have been blurred. We ask annotators to classify images as either “all faces blurred”
301 or “some unblurred faces”. To reduce the risk of automation bias, the gold standard checks for the
302 workers include large numbers of images for which one or more face detections have been dropped.
303 Thus, workers see a fairly large number of images with unblurred faces.

304 **Step 3 – Manual Spotting.** Many of the missing faces are simply faces from unusual angles that
305 are easily spotted by a human but understandably missed by a computer system. We ask workers to
306 annotate these with a box, focusing on faces that are clearly visible and large enough to recognize
307 (e.g., not 2 pixel tall faces in crowds).

308 **Step 4 – Verification of Manually Annotated.** We then run the images with the additional boxes
309 through our face blurring system, and ask workers to classify each image as either “all faces blurred”
310 or “some unblurred faces”. We apply similar gold standard checks to Step 2.

311 **Step 5 – Manual Annotation, Including Masks.** The remaining faces are difficult to annotate and
312 primarily depict outdoor scenes with many people. Many have a face or two missing in an otherwise
313 properly parsed scene. Some, however, show systematic failures where large numbers of people have
314 clearly visible faces that are large enough to be recognizable, but are entirely missed by the automatic
315 system. We hypothesize that these are due to systematic gaps in the training data.

316 These images are often complex and so the authors of the paper marked these images themselves.
317 Using photo editing programs, we marked regions to: (1) provide a bounding box for a face; (2)
318 provide a region that needed to be blurred. The ability to blur an entire region enables us to quickly,
319 for instance, blur a few hundred faces in tennis stands.

320 The boxes for redaction come from the union of Steps 1, 3, and 5. The final additional redaction
321 mask comes from Step 5.

322 Redaction Algorithm

323 We follow the strategy of [22], but make a few changes that catch some edge cases we observed in
324 our data. In particular, the photos we interact with often have fairly large ranges of depths of faces.

325 **Input.** As input, we are given an image $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$, a set of N boxes $B = \{(x_0^i, y_0^i, x_1^i, y_1^i)\}_{i=1}^N$
326 for the faces, and a redaction mask $\mathbf{R} \in [0, 1]^{H \times W}$ of pixels that will always be redacted.

327 **Data Prep.** First, we calculate a maximum diagonal across the boxes

$$d = \max_i \sqrt{((y_1^i - y_0^i)^2 + (x_1^i - x_0^i)^2)} \quad (1)$$

328 that defines the scale of the blur. This is used to calculate a universal blur filter for the entire image.

329 We then expand the boxes by a constant $c = 0.15d$, or $B' = \{(x_0^i - c, y_0^i - c, x_1^i + c, y_1^i + c)\}_{i=1}^N$.
330 These new boxes B' define the region that will be redacted.

331 We then create a pixelwise “not a face” mask $\mathbf{M} \in [0, 1]^{H \times W}$. $\mathbf{M}[x, y]$ is 1 if only the two things
332 are true: (1) (x, y) is outside a box in B' (i.e., was not marked as a box); (2) (x, y) is **not** marked as
333 to-redact in \mathbf{R} (i.e., was not marked as a redaction region).

334 **Blurring.** If we then define \mathbf{G} as a Gaussian filter with standard deviation $0.1d$, we compute an
335 alpha mask $\mathbf{A} = \mathbf{M} * \mathbf{G} \odot (1 - \mathbf{R})$, which smoothly blends from blurred to unblurred while also
336 hard-forcing anything inside a redaction mask to be blurred. Then the final image is

$$\mathbf{A} \odot \mathbf{I} + (1 - \mathbf{A}) \odot (\mathbf{I} * \mathbf{G}), \quad (2)$$

337 which uses $\mathbf{I} * \mathbf{G}$ (the blurred image) inside boxes and \mathbf{I} outside, with a smooth tradeoff (ex-
338 cept for redaction masks, where the cutoff is sudden). When we re-save images, we use PIL’s
339 `.save(..., quality="keep")` ability to re-use the DCT coefficients to avoid double-JPG artifacts.

340 Our algorithm differs slightly from [22] by expanding **all** boxes by a fraction of d . We found that
341 when faces varied in sizes, sometimes the redaction mask would get too blurred in $\mathbf{M} * \mathbf{G}$, and so
342 high frequency details of far away faces would peer through.

343 E.2 Child Detection

344 To mitigate concerns about the use of images containing children, even in publicly available and
345 creative-commons data, we filtered the data to remove children. We asked workers to annotate
346 whether any people under the age of 18 were in the picture. We only include images where annotators
347 came to a consensus that no children were in the picture. Both images of children and images with
348 inconsistent annotations were not included.

349 During the process of model development and data processing, we periodically came across pictures
350 of children (often in large crowds); we added these to the removed list.

351 We initially experimented unsuccessfully with an automatic approach that used face detection and age
352 regression. In short: we estimated the ages of people in the images based on detected faces and then
353 removed any face that was detected as a minor. We found this approach to be too inaccurate in terms
354 of both false positives and negatives. False positives (adults detected as children) were somewhat
355 idiosyncratic. False negatives were primarily people who were obviously children due to context
356 (e.g., clothing and size) but whose face were not visible or not clear.

357 F New Videos

358 We gathered a new video dataset using a semi-automatic method. This approach combines a small
359 amount of annotation with automatic approaches for relevance feedback. All the specifics about the
360 annotation instructions for this task appear in Section H.

361 F.1 Video Selection

362 We start with a collection of 9623 search terms generated combinatorially from Section F.3. For each
363 term, we return up to 200 videos (4 pages of results with up to 50 videos per page). **We search only**
364 **for videos explicitly marked as Creative Commons.** This returns 508,716 videos. We follow the
365 approach of [7, 18] where we use YouTube thumbnails to identify videos that are likely of interest.
366 The advantage of thumbnails is that they are substantially smaller than the video (typically under
367 100KB).

368 **Video representation for deep networks.** Given four thumbnail images from a video, we represent
369 each thumbnail with the final feature activation of an imagenet pre-trained resnet 50. We then
370 represent the video with an aggregation of the thumbnail representations. The video representation
371 is the concatenation of: the mean across dimensions, L2-normalized; the standard deviation across
372 dimensions; and the minimum, mean, and maximum of the pairwise distances between the feature
373 vectors.

374 We annotate two tasks in order to filter the videos.

375 **Task 1 Video Validity.** The annotator identifies unacceptable videos. Each annotator sees the
376 thumbnails montaged in a 2×2 grid. The annotator categorizes the video into one of three categories.

- 377 • (*Not Real*): These include cartoons, animations, screen recordings, slideshows, and
378 videogames. One or two thumbnails showing a diagram or logo (e.g., a subscription
379 request) is acceptable; more than this makes a video fall into “Not Real”.
- 380 • (*Lecturing*): These show a person sitting in front of the camera. If two or more frames show
381 the same person, in more or less the same posture, with the same background and talking to
382 the camera, then this is a “Lecturing” video.
- 383 • (*Shows Children*): If any of the thumbnails depict people who appear to be under 18, then
384 the video is classified as “Shows Children”.
- 385 • (*Acceptable*): Anything *other* than the above is considered acceptable.

386 We obtain 9,856 conclusively labeled samples, of which 6570 (66.7%) are “Acceptable”, 1824
387 (18.5%) are “Not Real”, 1169 (11.8%) are “Lecturing”, and 293 (3.0%) are “Shows Children”.

388 **Video Interaction.** The annotator identifies whether the video has at least two frames of hand-object
389 interaction. For each video, we extract its frames, and then make a 3×3 montage showing the frame
390 at 20%, 27.5%, 35%, 42.5%, 50%, 57.5%, 65%, 72.5%, and 80% of the way through the video. Note
391 that while the annotator sees nine frames (to see into the video to count), the network itself only has
392 access to the four thumbnails. The idea is that the network can learn correlations between how the
393 thumbnail is presented and the content in the video. Each annotator is asked to categorize:

- 394 • (*Interaction Rich*): If there are two or more frames that show a hand clearly engaged in
395 interaction (any form of contact other than resting on a table), then the video falls into this
396 category.
- 397 • (*Not Interaction Rich*): Any video showing fewer than two frames falls into this category,
398 including videos with no hands visible

399 We obtain 6,082 conclusively labeled samples, of which 4606 (75.7%) are “Interaction Rich” and
400 1476 (24.3%) are “Not Interaction Rich”.

401 **Filtering.** We fit two linear logistic regression model on the features. One predicts *Acceptable-vs-*
402 (either *Not Real* or *Shows Children*); the other predicts *Interaction Rich-vs-Not Interaction Rich*. We
403 then take $\approx 15,000$ videos from the intersection of the top 20% of the videos sorted by $p(\text{Acceptable})$
404 and the top 20% of the videos sorted by $p(\text{Interaction-Rich})$. We take random samples from the top
405 instead of the top predicted to ensure that our videos are representative of “interaction-rich” videos as
406 opposed to videos that maximally represent “interaction-rich”.

407 F.2 Frame Selection

408 Once we have selected $\approx 15,000$ videos, we extract one frame per second per video to generate a pool
409 of potential frames.

410 **Frame Representation.** We represent each frame using the final feature activation of an Imagenet
411 pre-trained Resnet-50.

412 **Scene Depth.** The annotator identifies the scene depth, split into three categories:

- 413 • (*Up Close*): This frame is probably within 50cm of the camera.
- 414 • (*Further*): This frame is probably at least 1m away. If hands are visible, they are at least 1m
415 away.
- 416 • (*Not Real Video*): This frame does not show a real frame (e.g., a diagram or text). We
417 provide this as an option to ensure consensus on the handful of non-real frames that are left.

418 Annotators are instructed to make a best guess for videos showing scenes with depths between 50cm
419 and 1m. All qualifiers and gold-standard tests are clearly in one category or another. We obtain 4979
420 conclusive samples, of which 3574 (71.7%) are *Up Close*, 1364 (27.4%) are *Further*, and 179 (3.6%)
421 are *Not Real Video*. We fit a multinomial logistic regression model to classify each video into these
422 categories. We then sample 50,000 frames randomly from the frames where $p(\text{Up Close}) > \frac{1}{2}$ and
423 50,000 frames randomly from the frames where $p(\text{Further}) > \frac{1}{2}$.

424 **F.3 Search Grammar**

425 We followed the following search grammar, following [18]. The data for each of the 12 categories is
426 generated by selecting a word from row, where ϵ is the empty string. Therefore, the DIY grammar
427 includes the searches “DIY IKEA genius“ and “furniture amazing“ and “creator hacks“.

428 Beauty:

- 429 • beauty, haircare, bodycare, make up, skincare
- 430 • routine, tips, tutorial, with me, secrets, ϵ
- 431 • morning, night, anti-aging, essential, affordable, at home, everyday, natural, realistic, ultimate, winter,
- 432 • summer, fall, autumn, spring, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, ϵ

433 Board Games:

- 434 • play, how to play, learn to play, win in
- 435 • board game, backgammon, checkers, chinese checkers, chess, darts, Go, halma, lotto, ludo, mah jongg,
- 436 • monopoly, pachisi, scrabble, shovel board, snakes and ladders, tic tac toe, tic-tak-toe
- 437 • 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, game, basic, beginners, master, guide, strategy

438 DIY:

- 439 • DIY, ϵ
- 440 • IKEA, gift, furniture, crafts, room, food, drink, decor, experiment, bag, waste, card, candy, cookie,
- 441 • desk, creator, boxes
- 442 • ideas, cheap, genius, master, amazing, office, home, random, hacks, 2015, 2016, 2017, 2018, 2019,
- 443 • 2020, 2021, 2022, ϵ

444 Drinks:

- 445 • made, make, kitchen, home made
- 446 • sip, tea, gulp, fizz, mate, milk, gulp, draft, cider, cocoa, mixer, coffee, cooler, posset, drinks, frappe,
- 447 • hydromel, smoothie, syllabub, wish-wash, refresher, ice milk, milkshake, soft drink, water, espresso,
- 448 • cappuccino, latte
- 449 • 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, ϵ

450 Food:

- 451 • make, cook, cooked, restaurant, home made
- 452 • meat, comfort food, pasta, bread, yolk, chocolate, foodstuff, baked goods, junk food, loaf, seafood,
- 453 • beverage, slop, fare, butter, comestible, produce, leftovers, miraculous food, soul food, feed, coconut,
- 454 • fish, food, yogurt, breakfast food, pizza, convenience food, cheese
- 455 • kitchen, restaurant, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, ϵ

456 Furniture:

- 457 • install, assembly, home
- 458 • nest, lamp, seat, table, buffet, cabinet, bedstead, etagere, washstand, bookcase, furniture, sectional,
- 459 • lawn furniture, chest of drawers, bedroom furniture, dining room furniture, wardrobe, ϵ
- 460 • home, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, ϵ

461 Gardening:

- 462 • backyard, indoor, garden, gardening, plant, grow
- 463 • care, vegetable, flower, tree, veggie, food, seed, greens, ϵ
- 464 • tips, idea, guide, spring, summer, fall, autumn, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, ϵ

465 Housework:

- 466 • clean, redo, housework, reorganize, decorate, tidy
- 467 • room, home bedroom, house, living room, dining room, apartment, home, ϵ
- 468 • motivation, tips, extreme, with me, routine, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, ϵ

469 Packing:

- 470 • pack, packing, unpack, unpacking, wrap, unbox
 - 471 • clothes, luggage, suitcase, bag, gift, lunch, food, travel, box, package, trip, cruise, vacation
 - 472 • essential, guide, tips, tricks, work, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, ϵ
- 473 Puzzles:
- 474 • solve, play, do
 - 475 • puzzle, jigsaw puzzle, sliding puzzle, jack puzzle, burr puzzle, lock puzzle, pyramid puzzle, ring
 - 476 puzzle, nail puzzle, lego, magic cube, Rubik's cube
 - 477 • beginner, impossible, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, ϵ
- 478 Repair
- 479 • repair, fix, maintain, maintenance
 - 480 • automobile, car, machine, trunk, mechanics, Jeep, vehicle, Ford, BMW, alternator, engine, bike,
 - 481 motorcycle, motor, generator, computer, pc, equipment, phone, earphone, watch, bulb, eletrics, electric
 - 482 appliance, l
 - 483 • 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022
- 484 Study:
- 485 • study, revise
 - 486 • with me, ϵ
 - 487 • exam, finals, midterms, midtest, dissertation, engineering, physics, history, psychology, economics,
 - 488 exam, finals, university
 - 489 • 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022

490 **G Datasheet for Hands23**

491 **Preamble**

492 The Hands23 dataset consists of annotations on four separate data sources: (1) a New Video dataset,
493 referred to as New Videos; (2) the 2017 training set for COCO [11]; (3) the frames from the Internet
494 Articulation Videos [16]; and (4) the training and validation frames of the EPIC-KITCHENS [5]
495 VISOR [6] challenge. Answering some of the standard datasheet questions involves answering
496 questions not just about the annotation, but also about the underlying data. Where it is relevant, we
497 have answered the question about the underlying data as well. The answers will be as follows:

498 **A.** This is an answer for the dataset

499 *A for New Videos.* This is an answer for the New Videos subset

500 *A for COCO.* This is an answer for the COCO subset

501 *A for Articulation.* This is an answer for the articulation subset

502 *A for VISOR.* This is an answer for the VISOR subset

503 **Motivation**

504 **Q. For what purpose was the dataset created? Was there a specific task in mind? Was there a**
505 *specific gap that needed to be filled? Please provide a description.*

506 **A.** This dataset, Hands23 was created to provide an improvement in the scale and quality of available
507 datasets for understanding hands interacting with the world. Past datasets have limitations in terms
508 of the richness of their annotation. As an ancillary benefit, many past datasets have included data
509 that was available under unclear copyright licenses and have included minors and unblurred faces.
510 Hands23 consists entirely of creative commons videos, blurs faces, and removes minors from the
511 data.

512 **Q. Who created the dataset (e.g., which team, research group) and on behalf of which entity**
513 **(e.g., company, institution, organization)?**

514 **A.** Cannot be answered during anonymous review but will be provided upon publication.

515 **Q. Who funded the creation of the dataset? If there is an associated grant, please provide the**
516 *name of the grantor and the grant name and number.*

517 **A.** Cannot be answered during anonymous review but will be provided upon publication.

518 **Q. Any other comments?**

519 **A.** No

520 **Composition**

521 **Q. What do the instances that comprise the dataset represent (e.g., documents, photos, people,**
522 **countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and**
523 *interactions between them; nodes and edges)? Please provide a description.*

524 **A.** The dataset contains images with corresponding annotations. The instances are hands and
525 associated information in these images. The corresponding paper *Towards A Richer 2D Understanding*
526 *of Hands at Scale* describes the annotations in more detail, but a brief description follows.

527 There are boxes for hands as well as additional annotation in terms of side, contact state, and (for
528 some hands) grasp type. Hands that are in contact with objects have a box for the in-contact box;
529 objects that are labeled as tools in use have a box for the object that tool is in contact with.

530 **Q. How many instances are there in total (of each type, if appropriate)?**

531 **A.** There are approximately 257K images containing annotations of 400K hands, 288K objects, and
532 19K second objects.

533 **Q. Does the dataset contain all possible instances or is it a sample (not necessarily random)**
534 **of instances from a larger set?** *If the dataset is a sample, then what is the larger set? Is the*
535 *sample representative of the larger set (e.g., geographic coverage)? If so, please describe how*
536 *this representativeness was validated/verified. If it is not representative of the larger set, please*
537 *describe why not (e.g., to cover a more diverse range of instances, because instances were withheld*
538 *or unavailable).*

539 **A.** There are a few downsampling steps in the creation of the dataset. For annotations, the only
540 downsampling done is in not labeling all hands with grasp information and not labeling far-away
541 hands in COCO. Grasp annotation is expensive, and so a random subset of hands were annotated
542 with grasps. Far-away hands are hard to see, and so were not annotated in COCO (specifically: only
543 non-crowd ≥ 1000 pixel persons had their hands annotated) The important downsampling happened
544 with in image selection. We report what we know about each dataset below, but note that people who
545 appear to be minors have been removed from all data.

546 *A for New Videos.* The data was selected from a large collection of videos described in the supplement
547 for the paper. Once videos were selected, frames were selected randomly subject, subject to an
548 automated balancing of estimated overall scene depth.

549 *A for Articulation.* Articulation data comes from videos that appears to have been selected using a
550 procedure that appears to be similar to New Videos according to [16].

551 *A for VISOR.* According to its documentation, VISOR frames were sampled to be denser within
552 actions and then further selected to have reduced blur.

553 *A for COCO.* COCO data was gathered using the COCO pipeline.

554 **Q. What data does each instance consist of?** *“Raw” data (e.g., unprocessed text or images) or*
555 *features? In either case, please provide a description.*

556 **A.** The dataset consists of images with annotations. The core instance for the dataset is a hand in an
557 image. This hand has:

- 558 1. a box location;
- 559 2. left-vs-right as a binary classification;
- 560 3. contact state as a multi-class classification into: {*no contact, self-contact, other-contact,*
- 561 *object contact*};
- 562 4. fine-grained contact state as a multi-class classification into: { *tool-used, tool-held, tool-*
- 563 *touched, container-held, container-touched, neither-held, neither-touched* };
- 564 5. a box for the contacted object if the hand is in contact;
- 565 6. a box for the object that the contacted object is in contact with if the contacted object is a
- 566 tool;
- 567 7. grasp information, for a random subset of hands in contact with objects, which is framed as
- 568 a multi-class classification problem into { *Non-Prehensile-Fingers-Only, Non-Prehensile-*
- 569 *Palm, Power-Prismatic, Power-Circular, Precision-Prismatic, Precision-Circular, Lateral*}.

570 Every image has zero or more hands with these annotations.

571 **Q. Is there a label or target associated with each instance?** *If so, please provide a description.*

572 **A.** Yes. Please see the above.

573 **Q. Is any information missing from individual instances?** *If so, please provide a description,*
574 *explaining why this information is missing (e.g., because it was unavailable). This does not include*
575 *intentionally removed information, but might include, e.g., redacted text.*

576 **A.** Yes, for two reasons. First, only a subset of hands were labeled with grasps because grasp
577 annotation is difficult and expensive. Second, annotators could not come to a consensus on some
578 annotations. These are marked in the dataset as unknown.

579 **Q. Are relationships between individual instances made explicit (e.g., users’ movie ratings,**
580 **social network links)?** *If so, please describe how these relationships are made explicit.*

581 **A.** The hand instances are not linked to each other, but the object instances are linked to the hands,
582 and the second objects instances are linked to the objects. The objects were linked to the hands
583 explicitly through the annotation process: the objects are labeled as “what is the object that is in
584 contact with this hand”.

585 **Q. Are there recommended data splits (e.g., training, development/validation, testing)?** *If so,*
586 *please provide a description of these splits, explaining the rationale behind them*

587 **A.** Yes, we provide recommended data splits that are chosen carefully. Our splits are chosen to
588 minimize the chance of source contamination (e.g., data from the same channel appearing in multiple
589 splits) and maximize the agreement with existing datasets.

590 We split the source datasets as follows:

591 *A for COCO.* We annotate the training set of COCO 2017. We assign COCO images randomly into
592 our training set (80%) of images, validation set (10% of images), and test set (10% of images).

593 *A for VISOR.* VISOR has a held-out test set that we do not annotate. We follow the VISOR split as
594 follows: we make the VISOR validation set our test set; we then randomly assign VISOR’s training
595 set into our training set (80% of images) and our validation set (20% of images)

596 *A for Articulation.* We try to follow the split in [16] as closely as possible. However, if we know that
597 two frames come from the same channel, we assign them to the same split. The split promotion logic
598 is: if the channel contains a test frame, then all the channel’s frames are moved to test; if a channel
599 contains no test frames and at least one validation frame, then all the channel’s frames are moved to
600 validation.

601 *A for New Videos.* We split the videos by channel, aiming to assign 80% to train, 10% to validation,
602 and 10% to test. In other words, all of a channel’s frames are in only one split. We assign channels
603 to the split randomly, except for videos that appear in the Articulation dataset (which are assigned
604 according to the Articulation splits).

605 **Q. Are there any errors, sources of noise, or redundancies in the dataset?** *If so, please provide a*
606 *description.*

607 **A.** There are likely incorrect annotations in the dataset, as is the case with all annotations. There are
608 no deliberate redundancies beyond what is present when annotating frames from videos.

609 **Q. Is the dataset self-contained, or does it link to or otherwise rely on external resources**
610 **(e.g., websites, tweets, other datasets)?** *If it links to or relies on external resources, a) are there*
611 *guarantees that they will exist, and remain constant, over time; b) are there official archival versions*
612 *of the complete dataset (i.e., including the external resources as they existed at the time the dataset*
613 *was created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external*
614 *resources that might apply to a future user? Please provide descriptions of all external resources and*
615 *any restrictions associated with them, as well as links or other access points, as appropriate.*

616 **A.** The dataset depends on a few different source datasets. We will provide an archival purpose of the
617 dataset for non-commercial research purposes. We will not charge a fee, but users must agree to the
618 restrictions of the underlying data. The link for download is not available at the time of submission
619 of the paper.

620 **Q. Does the dataset contain data that might be considered confidential (e.g., data that is**
621 **protected by legal privilege or by doctor–patient confidentiality, data that includes the content**
622 **of individuals’ non-public communications)?** *If so, please provide a description.*

623 **A.** We do not believe so and did not find any during our use of the dataset. For Internet data, the data
624 was posted publicly by users on photo/video sharing websites, and we expect that users would have
625 exercised precaution. For VISOR, the capture process involves the user watching and verifying their
626 own data, so we expect that users would have also exercised caution.

627 **Q. Does the dataset contain data that, if viewed directly, might be offensive, insulting, threaten-**
628 **ing, or might otherwise cause anxiety?** *If so, please describe why*

629 **A.** We do not believe so. However, what causes anxiety will differ from person to person – if people
630 find videos of cooking animal meat anxiety-incuding, for instance, there are videos of this in the
631 dataset.

632 **Q. Does the dataset relate to people?** *If not, you may skip the remaining questions in this section.*

633 **A.** Yes. The dataset relates to people.

634 **Q. Does the dataset identify any subpopulations (e.g., by age, gender)?** *If so, please describe how*
635 *these subpopulations are identified and provide a description of their respective distributions within*
636 *the dataset.*

637 **A.** No. We do not identify demographic information of people in the dataset, except for a post-hoc
638 audit of model performance.

639 **Q. Is it possible to identify individuals (i.e., one or more natural persons), either directly or**
640 **indirectly (i.e., in combination with other data) from the dataset?** *If so, please describe how*

641 **A.** Although we have taken the steps to obfuscate faces in the data, it is certainly possible for a person
642 with time to identify users from the data. First, the data license for our data is creative commons,
643 which requires attribution. This attribution intrinsically may help identify users. Second, the data
644 itself was public on a video sharing website, so we are not releasing new data. However, we believe
645 that the face obfuscation and removal of minors from the data provides some privacy.

646 **Q. Does the dataset contain data that might be considered sensitive in any way (e.g., data**
647 **that reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or**
648 **union memberships, or locations; financial or health data; biometric or genetic data; forms**
649 **of government identification, such as social security numbers; criminal history)?** *If so, please*
650 *provide a description.*

651 **A.** It is possible that some information can be gleaned from the videos. However, this is data that
652 users had uploaded and therefore the amount of information that is given away is not more than what
653 previously was published to the Internet.

654 **Collection Process**

655 **Q. How was the data associated with each instance acquired?** *Was the data directly observ-*
656 *able (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly*
657 *inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or lan-*
658 *guage)? If data was reported by subjects or indirectly inferred/derived from other data, was the data*
659 *validated/verified? If so, please describe how.*

660 **A.** The data is a combination of images (which were directly obtained from Internet data or existing
661 datasets) as well as annotated labels. The labels were annotated by multiple workers using standard
662 labeling protocols (a qualifier to verify task understanding, checks to verify correct annotations, and
663 multiple judgments to check for annotation consensus). The resulting annotations were checked for
664 correctness during the process by researchers on the project.

665 **Q. What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or**
666 **sensor, manual human curation, software program, software API)?** *How were these mechanisms*
667 *or procedures validated?*

668 **A.** The annotations were obtained primarily by working with a crowdsourcing company. The precise
669 process is documented in the supplemental materials. The images themselves were obtained as
670 follows.

671 *A for New Videos.* The data was obtained with custom scripts for scanning for Creative Commons
672 videos on YouTube.

673 *A for Articulation.* Unknown to us and not listed by the authors; we expect it is similar to New Videos.

674 *A for VISOR.* Recorded with collaboration of the depicted people, according to the datasheet for
675 VISOR.

676 *A for COCO.* Unknown to us and not listed by the authors.

677 **Q. If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,**
678 **probabilistic with specific sampling probabilities)?**

679 **A.** The only subsampling done in annotations is in subsampling which grasps were annotated. This
680 was done at random. The images themselves were subsampled, which we report in the question on
681 Composition.

682 **Q. Who was involved in the data collection process (e.g., students, crowdworkers, contractors)**
683 **and how were they compensated (e.g., how much were crowdworkers paid)?**

684 **A.** Data collection involved both the researchers and crowdworkers that a third party company hired.
685 *Researchers.* Researchers involved in the project did pilot annotations of data and the final face
686 blurring efforts.

687 *Crowdworkers.* We hired a third party company to annotate the data. This company performs
688 annotation of a set of discrete tasks (e.g., categorization, boxes, segmentation). The use of a third
689 party intermediary makes it hard to estimate compensation, but for transparency, we report the
690 annotation budget and breakdown into categories in the supplementary material of the paper.

691 **Q. Over what timeframe was the data collected? Does this timeframe match the creation timeframe**
692 **of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please**
693 **describe the timeframe in which the data associated with the instances was created.**

694 **A.** The annotations were collected over a nearly year-long period from late June 2022 until early May
695 2023. The individual data for each dataset was collected:

696 *A for New Videos.* The data was scanned and downloaded early-to-mid-September 2022.

697 *A for Articulation.* Unknown and not listed by the authors.

698 *A for VISOR.* April 2017 through July 2020, according to the datasheet for VISOR.

699 *A for COCO.* Not listed, but we presume no later 2017 and likely close to 2017.

700 For VISOR, the collection timeframe matches the creation timeframe; for others, the collection
701 timeframe does not match the creation timeframe. The timestamps, for instance, on the video
702 downloads for New Videos suggest that some videos may be as old as 2010. Judging by the image
703 content in COCO, this data was likely captured far before 2017.

704 **Q. Were any ethical review processes conducted (e.g., by an institutional review board)? If so,**
705 **please provide a description of these review processes, including the outcomes, as well as a link or**
706 **other access point to any supporting documentation**

707 **A.** For Internet data, there were no formal review processes followed because the data was pre-existing
708 and public and did not involve interaction with the participants. VISOR is based on EPIC-KITCHENS,
709 which involved interaction with participants. EPIC-KITCHENS was collected with University of
710 Bristol faculty ethics approval. These application is held at the university of Bristol and the participant
711 consent form is available here

712 **Q. Does the dataset relate to people? If not, you may skip the remaining questions in this section.**

713 **A.** Yes. The dataset contains people.

714 **Q. Did you collect the data from the individuals in question directly, or obtain it via third parties**
715 **or other sources (e.g., websites)**

716 **A.** This depends on the source of data. New Videos and Internet Articulation [16] come from YouTube
717 via searching for CreativeCommons-licensed data. COCO comes from similarly searching Flickr.com
718 for CreativeCommons-licensed data. VISOR was collected directly by and with the individuals
719 depicted.

720 **Q. Were the individuals in question notified about the data collection? If so, please describe (or**
721 **show with screenshots or other information) how notice was provided, and provide a link or other**
722 **access point to, or otherwise reproduce, the exact language of the notification itself.**

723 **A.** No for the Internet data (New Videos, Internet Articulation, COCO); yes for VISOR.

724 In the case of Internet data, users had posted this data publicly to websites meant for sharing photos
725 and videos and selected a CreativeCommons license. Thus the users who captured the photos
726 presumably knew that their data would be public, but were not explicitly informed that their data
727 would be used for machine learning research. As a mitigation for concerns about data use, we remove
728 minors from the dataset and blur all the faces.

729 In the case of VISOR, yes. Since the data was directly collected by the participants, the participants
730 were aware of the data collection process. All participants were given the opportunity to ask questions
731 before participating, and they could withdraw at any time without giving a reason. Participants
732 consented to the process and watched their footage. All participants were volunteers and were not
733 compensated.

734 **Q. Did the individuals in question consent to the collection and use of their data?** *If so, please*
735 *describe (or show with screenshots or other information) how consent was requested and provided,*
736 *and provide a link or other access point to, or otherwise reproduce, the exact language to which the*
737 *individuals consented.*

738 **A.** Similar to the above: for Internet data, no consent was obtained but the data was previously made
739 public and we have removed minors and blurred faces. For VISOR, the participants consented to data
740 collection and use and reviewed their footage before its use.

741 **Q. If consent was obtained, were the consenting individuals provided with a mechanism to**
742 **revoke their consent in the future or for certain uses?** *If so, please provide a description, as well*
743 *as a link or other access point to the mechanism (if appropriate).*

744 **A.** For Internet data, we will provide a mechanism to remove data from the dataset for users upon
745 release of the data. For VISOR: participants were able to withdraw from the process at any point
746 until the data was published by DOI. At the moment, participants are unable to withdraw their data.

747 **Q. Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data**
748 **protection impact analysis) been conducted?** *If so, please provide a description of this analysis,*
749 *including the outcomes, as well as a link or other access point to any supporting documentation.*

750 **A.** For Internet data, no. For VISOR, the University of Bristol faculty ethics committee reviewed the
751 protocol, and approved the dataset. They checked any potential impact and as the data is anonymous
752 no further actions were deemed as needed.

753 **Q. Any other comments?**

754 **A.** No

755 **Preprocessing/cleaning/labeling**

756 **Q. Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,**
757 **tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing**
758 **of missing values)?** *If so, please provide a description. If not, you may skip the remainder of the*
759 *questions in this section.*

760 **A.** We blurred faces in all of the data except for VISOR (which has no faces). The data is otherwise
761 unchanged (apart from basic format processing steps such as converting videos to frames).

762 **Q. Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to**
763 **support unanticipated future uses)?** *If so, please provide a link or other access point to the “raw”*
764 *data.*

765 **A.** The only raw data that exists before our images are: (a) the original source videos for the video
766 datasets; and (b) the frames with unblurred faces. We do not plan to publicly release the frames with
767 unblurred faces.

768 **Q. Is the software used to preprocess/clean/label the instances available?** *If so, please provide a*
769 *link or other access point.*

770 **A.** Not at the moment. Most of the software is one-off scripts that are not likely not of interest due to
771 their simplicity and non-general purpose nature. However, we are happy to share the code used for
772 blurring upon request.

773 **Q. Any other comments?**

774 **A.** No

775 **Uses**

776 **Q. Has the dataset been used for any tasks already?** *If so, please provide a description.*

777 **A.** Yes, the dataset has been used for hand detection, as documented in the paper. This task requires
778 localizing hands, the objects they hold, and the objects that are being touched by tools they use.
779 Additionally, the task requires predicting a variety of extra properties such as contact state and grasp
780 type.

781 **Q. Is there a repository that links to any or all papers or systems that use the dataset?** *If so,
782 please provide a link or other access point*

783 **A.** No, not at the moment.

784 **Q. What (other) tasks could the dataset be used for?**

785 **A.** We envision that the dataset could be used for a wide variety of other tasks. Earlier datasets in this
786 area have been used for tasks such as unsupervised learning for robotics.

787 **Q. Is there anything about the composition of the dataset or the way it was collected and
788 preprocessed/cleaned/labeled that might impact future uses?** *For example, is there anything that
789 a future user might need to know to avoid uses that could result in unfair treatment of individuals or
790 groups (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial harms,
791 legal risks) If so, please provide a description. Is there anything a future user could do to mitigate
792 these undesirable harms?*

793 **A.** There are a number of important considerations for using this dataset.

794 First, the dataset is not necessarily representative of the world's demographics: due to the collection
795 process, our data primarily reflects the users of YouTube and Flickr, and our egocentric data mainly
796 comes from the EPIC-KITCHENS benchmark. If the system is used in scenarios where accuracy is
797 critical, we would urge future users to do an evaluation on their data to make sure that there are no
798 biases in terms of performance.

799 Second, regardless of demographics, the dataset does not represent real-life due to the source of data.
800 Some of this lack of realism is missing data: COCO images rarely show transitional moments when
801 a tool is being used to interact with an object. Other lack of realism is due to realistic interactions
802 being chained together in unrealistic ways. For instance, New Videos contains many videos of people
803 attempting to eat enormous amounts of food. The interaction of picking up a piece of pizza may be
804 realistic, but the number of slices of pizza may not be.

805 Finally, the released data and models are trained on data with blurred faces. We find that unblurred
806 faces are occasionally seen as hands. Future users may wish to either preemptively blur faces going
807 into the model, or suppress detections that overlap with faces.

808 **Q. Are there tasks for which the dataset should not be used?** *If so, please provide a description.*

809 **A.** VISOR requires non-commercial research use only and so the full dataset can only be used for
810 non-commercial purposes. A commercial license for VISOR can be acquired through negotiation
811 with the University of Bristol.

812 Additionally, given the unrealistic nature of some of the underlying data, we would caution drawing
813 conclusions from the dataset in terms of frequency of events or how people interact with objects.

814 **Q. Any other comments?**

815 **A.** No

816 **Distribution**

817 **Q. Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,**
818 **organization) on behalf of which the dataset was created? If so, please provide a description.**

819 **A.** Yes. The dataset will be available for non-commercial purposes publicly.

820 **Q. How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the**
821 **dataset have a digital object identifier (DOI)?**

822 **A.** The dataset will be released via the project website with a to-be-determined format. We will also
823 provide a DOI.

824 **Q. When will the dataset be distributed?**

825 **A.** Not known at this point

826 **Q. Will the dataset be distributed under a copyright or other intellectual property (IP) license,**
827 **and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and**
828 **provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU,**
829 **as well as any fees associated with these restrictions.**

830 **A.** The underlying data of VISOR requires this dataset to have a Creative Commons BY-NC 4.0
831 license, which restricts commercial use of the data.

832 **Q. Have any third parties imposed IP-based or other restrictions on the data associated with**
833 **the instances? If so, please describe these restrictions, and provide a link or other access point**
834 **to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these**
835 **restrictions.**

836 **A.** There are no restrictions from third parties on the dataset.

837 **Q. Do any export controls or other regulatory restrictions apply to the dataset or to individual**
838 **instances? If so, please describe these restrictions, and provide a link or other access point to, or**
839 **otherwise reproduce, any supporting documentation.**

840 **A.** No. There are no restrictions beyond following the underling licenses of the image datasets

841 **Q. Any other comments?**

842 **A.** No

843 **Maintenance**

844 **Q. Who will be supporting/hosting/maintaining the dataset?**

845 **A.** The dataset will be released via a scheme that enables long-term preservation of the data even if
846 there are personnel changes. The precise details cannot be revealed at the moment to preserve the
847 anonymity of the work.

848 **Q. How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

849 **A.** The creators of the dataset will be listed in the corresponding paper and can be contacted via email
850 once their identities are made public.

851 **Q. Is there an erratum? If so, please provide a link or other access point.**

852 **A.** Not yet. If there are errata or updates, we will provide them on the dataset website once released.

853 **Q. Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete**
854 **instances)? If so, please describe how often, by whom, and how updates will be communicated to**
855 **users (e.g., mailing list, GitHub)?**

856 **A.** We do not have concrete plans as of yet; we will announce any updates on the dataset website
857 once released.

858 **Q. If the dataset relates to people, are there applicable limits on the retention of the data**
859 **associated with the instances (e.g., were individuals in question told that their data would be**
860 **retained for a fixed period of time and then deleted)?** *If so, please describe these limits and explain*
861 *how they will be enforced*

862 **A.** There are no limits on the retention of data at this point. We will monitor best practices and
863 re-assess after one year.

864 **Q. Will older versions of the dataset continue to be supported/hosted/maintained?** *If so, please*
865 *describe how. If not, please describe how its obsolescence will be communicated to users.*

866 **A.** For some changes yes. If we provide corrections to annotations or other updates that are not
867 intended to be removing data, we will have version control. If we remove data (e.g., due to offensive
868 imagery discovered), we will not provide public access to older versions.

869 **Q. If others want to extend/augment/build on/contribute to the dataset, is there a mechanism**
870 **for them to do so?** *If so, please provide a description. Will these contributions be validated/verified?*
871 *If so, please describe how. If not, why not? Is there a process for communicating/distributing these*
872 *contributions to other users? If so, please provide a description.*

873 **A.** Users are free to extend the dataset on their own and create derivative works, so long as they follow
874 the license agreements of the data. There is no official mechanism to incorporate new contributions,
875 but we encourage others to email us to let us know about extensions and modifications

876 **Q. Any other comments?**

877 **A.** No

878 H Data Annotation and Instructions

879 We now describe the data annotation process. We used an annotation company, HIVE (known as
880 *thehive.ai*) for nearly all annotation except some complex tasks that were done by the authors.

881 **Quality Control.** HIVE implements standard quality control mechanisms during the annotation
882 process. These consist of qualifiers (tests during the instructions to ensure that the instructions
883 are understood), gold standard checks (tests during annotation to ensure annotation quality), and
884 consensus labeling of judgments. Gold standard checks were selected specifically to be non-tricky
885 judgments: the goal was to serve as a sentinel to catch random guessing.

886 While the platform does not permit two-way interaction with annotators, we carefully monitored the
887 annotators to identify if our instructions were unclear, tasks were unfair, or if there were other issues.
888 We did this by monitoring performance on qualifier and gold-standard checks to find and remove
889 ambiguous annotations. We also reviewed the free-form feedback and ratings of our tasks by the
890 annotators. These free-form annotations often described the clarity of the instructions, whether they
891 thought that compensation was in line with their expectations, and difficulty.

892 **Compensation and Annotator Backgrounds.** The overall annotation budget of the project was
893 approximately \$40,000. Due to the use of an intermediary, converting our spent dollars into hourly
894 rates is difficult. However, in this section, we aim to provide as much transparency about how much
895 was spent and on what, including detailed information about the cost of each subtask.

896 Due to the nature of the platform, we do not have information the location or demographics of the
897 annotators. However, given that our tasks are primarily questions that are concretely defined in terms
898 of physical properties, we do not expect that annotator demographics will have a large impact.

899 **Instruction Screenshots.** We include annotation instruction screenshots. These have also had
900 their faces blurred and children redacted for consistency with the paper, but the annotators saw the
901 unredacted picture.

902 Naturally, it is difficult to show instructions for spotting unblurred faces or spotting children when
903 the faces in this document are blurred and children are removed. When there are faces that were
904 *not blurred* to illustrate unblurred faces for annotators, they are indicated with a black and white
905 checkerboard pattern. When children are removed in this section, we hide it with a black mask. One
906 face is left unredacted to illustrate the instructions. This face belongs to Nicolas Cage, who is a
907 celebrity.

908 **H.1 Hand Detection**

909 Approximately \$8,400 was spent on 280K box labeling tasks for hand detection. The particular
910 strategy depended per dataset. We used one strategy for VISOR, another for Articulation and New
911 Videos, and finally a third strategy for COCO. In both cases, annotators marked boxes and indicated
912 the side (left-vs-right) at the same time.

913 We examined the individual results and found that some hands were correct but had only one annotator
914 marking them. We later obtained contact state for these hands using a task that also had a not a hand
915 option. Hands marked with a valid contact state were kept; hands marked with “not a hand” were
916 discarded. This provided an additional set of boxes that let us reach our final number

917 **H.1.1 Annotation for VISOR**

918 For VISOR data, we simply ask workers to bound hands. There is little ambiguity.

Please draw bounding boxes around hands in the image, indicating left vs right for the hand.

We will show you two images stacked on top of each other.

- The top image will indicate which people we are interested in. Please do not annotate on the top image. We are providing it to make it easier to see the image when you annotate.
- The bottom image is where you should annotate. Please do annotate on the top image

Here's how to annotate:

- We are only interested in a some people in the image. We've outlined them in red and blue in the top. Please annotate hands **only** if they belong to a person that is outlined in the top of the image
- Please draw the tightest bounding box that includes all of the hand pixels, including the wrists
- If some parts of the hand are not visible, then please draw the **smallest** bounding box that encloses all visible parts of the hand. Please do not make any predictions about the non-visible parts.

919

920 **H.1.2 Annotation for Articulation, New Videos**

921 For egocentric data, we also reminded workers that left and right are mirrored depending on the
922 camera view

Draw draw bounding boxes around hands in the image, indicating left vs right.

- Please draw the tightest bounding box that includes all of the hand pixels, including the wrists
- If some parts of the hand are not visible, then please draw the **smallest** bounding box that encloses all visible parts of the hand.
- Please do not make any predictions about the non-visible parts.
- Please be careful about left vs right: some of the cameras are on the person's body (first person view) and some of the cameras are pointing towards the person's body (third person view)

Here are three examples.

Example 1:



923

Example 2: Note that this picture is taken from the point of the view of the person who's hand this is.



Example 3:



Example 4: Note that while the object hides part of the hand, you should annotate the full hand with a single box.



924 **H.1.3 Annotation for COCO**

925 For COCO, workers annotated only hands for humans who were non-crowd and had at least 1000
926 pixels in area. We indicated these humans with a red and blue border as shown below. Workers were
927 paid more for for images with more people; this rate was adjusted mid-project.

Please draw bounding boxes around hands in the image, indicating left vs right for the hand.

We will show you two images stacked on top of each other.

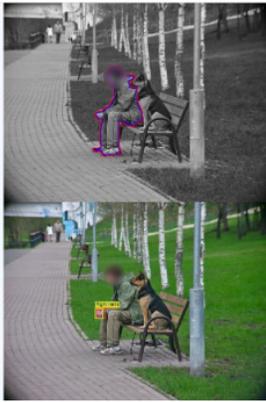
- The top image will indicate which people we are interested in. **Please do not annotate on the top image.** We are providing it to make it easier to see the image when you annotate.
- The bottom image is where you should annotate. **Please annotate on the bottom image.**

Here's how to annotate:

- We are only interested in a some people in the image. We've outlined them in red and blue in the top image. Please annotate hands **only** if they belong to a person that is outlined in the top of the image. Again, please annotate in the bottom image.
- Please draw the tightest bounding box that includes all of the hand pixels, including the wrists
- If some parts of the hand are not visible, then please draw the **smallest** bounding box that encloses all visible parts of the hand. **Please do not make any predictions about the non-visible parts.** The only exception to this rule are things like gloves and mittens where you know precisely where the hand is.

For example, in this below image:

- Please do not annotate the top image. This is just so we can indicate which people we would like annotated
- Please do not annotate the people who are far away. We have not outlined them in red and blue.



928

929 **H.2 Hand Contact State**

930 Approximately \$5,800 was spent on 500K hand contact tasks. This was framed as a standard
931 classification task between no contact, self-contact, other person contact, and object contact.

932 **H.2.1 VISOR, Articulation, New Videos**

933 We had relatively simple annotations for these examples.

1 2

Contact

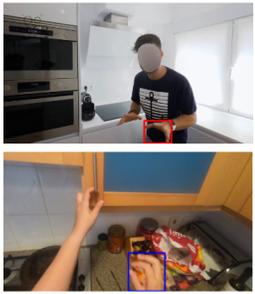
We'd like you to categorize the physical contact of a person's hands. By physical contact, we mean touching. We will draw a red or blue box around the hand that we are asking about. Left hands are drawn in red. Right hands are drawn in blue. This coloring may help if there are two hands near each other.

There are four categories:

- **No contact:** the hand is not in contact with anything. For example: a person is waving to a friend or pointing at something.
- **Self-contact:** the hand is contacting some other part of the body. For example: a person is washing their hands, or rubbing their neck.
- **Other person contact:** the hand is contacting someone else's body. For example: a person is shaking someone's hand, or hugging their friend.
- **Object contact:** the hand is contacting a non-human object. For example: a person is holding a knife, or opening a door.

We will next show you examples of these cases

No contact



Self-contact



Other person contact



Object Contact



935 **H.2.2 COCO**

936 We provided more examples for COCO.

1

Contact

We'd like you to categorize the physical contact of a person's hands. By physical contact, we mean touching. We will draw a red or blue box around the hand that we are asking about in the top image. Left hands are drawn in red. Right hands are drawn in blue. This coloring may help if there are two hands near each other. We will also show you the original image below with no box to help you better see the hand.

There are four categories:

- No contact: the hand is not in contact with anything. For example: a person is waving to a friend or pointing at something
- Self-contact: the hand is contacting some other part of the body. For example: a person is washing their hands, or rubbing their neck
- Other person contact: the hand is contacting someone else's body. For example: a person is shaking someone's hand, or hugging their friend
- Object contact: the hand is contacting a non-human object. For example: a person is holding a knife, or opening a door

We will next show you examples of these cases

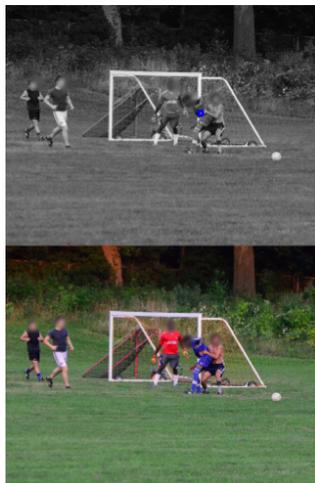
No contact



937

3

Other person contact



938 H.3 Additional Annotations – Checking Hands Labeled by One Annotator

939 We ran an additional annotation for hands that were annotated in the box detection stage by only one
 940 annotator. We asked for the contact information or for annotators to indicate that no hand was present.
 941 Hands that were annotated with a contact state were kept as hands. The information for the additional
 942 possible label is shown below:

2

Self-contact



4

Object Contact



Not a hand

This box is not around a hand. Please mark incorrect boxes as "not a hand".



943

953 **H.5 Object Tool/Container Status**

954 Approximately \$5,700 was spent obtaining tool and container status over 310K classification tasks.
 955 This was framed as a classification task.

956 **Summary.**

Overview
 Welcome!
 We are trying to classify the objects people are holding in videos (with the object colored in red) and how it's being held by a hand (with the object colored in blue). We are classifying the hand and the object that is being held according to two properties.

1. **Type: what type is the object?** There are three types: tools, containers, and neither (not a tool or container). Tools can be used to interact physically with something (ex. a knife, spoon) and containers can contain something (ex. or a bowl, bag, or mug). Objects like laptops, bricks, tables, chairs, and cell phones aren't tools or containers.
2. **Use: what is the object being used for?** There are three options: used, held/carried, and touched. All types of objects can be held/carried (ex. a knife, bowl, or brick in someone's hand) or touched (ex. someone with their hands on a knife, bowl, or brick). Tools can also be used, or be in active use (ex. a knife cutting a potato).

Together this gives seven categories: three for tools, two for containers, and two for neither.

When we are asking the question about the blue hand and the red object. For instance, in this example, the **container** is being held by the **hand in blue**, but just touched by the hand that's not in blue.



957

958 **Tools Section.**

1

Tools
 Tools are objects people can use to do things in the physical world. We are interested in physical tools that can be used to physically act like a hand and are not containers.

- Examples of tools: knives, spoons, screwdrivers, ladles, baseball bats, tennis rackets, paddles are tools.
- Examples of things that are clearly not tools: bricks, tables, beds, pomegranates, flowers.
- Examples of things that are not tools in our definitions: bowls (since it's a container), cell phones, light switches, game console controllers, cameras (since they don't physically enable you to use the tool like a hand), umbrellas, clothing, or hats (since they aren't used like a hand).

For the tools, there are three states we are interested in

1: Tool, Used
 If the tool is in contact with another object, it's in use. For instance, this spoon is in contact with the food in the pot, so it's a tool in use.

959



2

Tools
 Tools are objects people can use to do things in the physical world. We are interested in physical tools that can be used to physically act like a hand and are not containers.

- Examples of tools: knives, spoons, screwdrivers, ladles, baseball bats, tennis rackets, paddles are tools.
- Examples of things that are clearly not tools: bricks, tables, beds, pomegranates, flowers.
- Examples of things that are not tools in our definitions: bowls (since it's a container), cell phones, light switches, game console controllers, cameras (since they don't physically enable you to use the tool like a hand), umbrellas, clothing, or hats (since they aren't used like a hand).

For the tools, there are three states we are interested in

1: Tool, Used
 If the tool is in contact with another object, it's in use. For instance, this spoon is in contact with the food in the pot, so it's a tool in use.



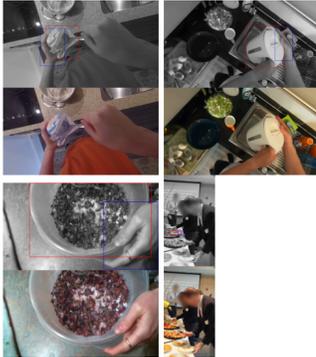
960 **Containers Section.**

1

Containers
 Containers are objects that can contain other objects. For instance, a bowl, plate, tray, bag, bottle and bin can contain other objects. For containers, there are also two states we are interested in: held vs not held. *Please note that the qualifier for this section will contain one tool.*

4. Container, Held/Carried
 If a container is held/carried it's held/carried. This is true regardless of whether it's in use or empty. For instance, all four are held/carried

961



2

5. Container, Touched
 If a container is not actually held, but just incidentally touched, then it's touched. For instance, the person isn't actually holding the container but is just moving it around by touching it.



962 **Neither Section.**

1

Neither Tool nor Container

Some objects can't be used as a tool or a container. For instance, bananas, bricks, tables, and chairs.

6: Neither, Held

If the non-tool/container object is held, then it's Neither, Held

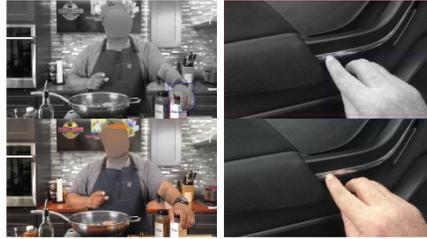
963



2

7: Neither, Touched

If the object is just touched, but not held, then it's neither, touched.



964 **H.6 Second Box**

965 Approximately \$1,200 was spent on obtaining boxes for a second object over 46K box tasks. This
 966 was done similarly to annotating the in-contact box: we provided the hand and object in the top half
 967 of the image and annotators annotated the bottom half of the image.

1

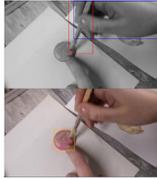
Bounding Objects Tools are interacting

We are looking to understand people interacting with objects using tools. In particular, we're trying to get the extent of objects that are being touched by tools. For each image, we will show in the top half of the image:

- a hand: this is in blue
- an object: this will be shown in red and is often a tool that someone is using to achieve a goal. This tool is probably touching another object

We would like you to draw a bounding box around the bottom of the image bounding the **interacted-with object**: this is the object that the tool is touching. For example, given the hand touching the paintbrush tool, we would like you to draw a box around the object in the bottom.

Please include only the visible parts of the object. Do not guess about parts of the object that you cannot see.

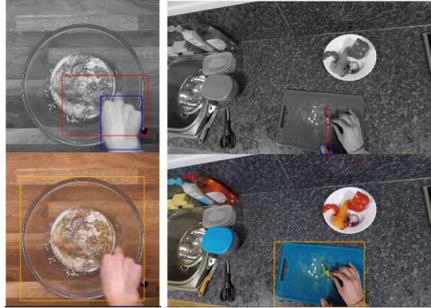


2

Special Instructions (important)

Containers and cutting boards:

- If the tool is inside a container like a bowl or pan, please include the full container.
- If the object is on a cutting board and it is not obvious whether the knife is contacting the cutting board or not, please include the cutting board.
- If the object is just sitting on a surface (such as in the first example) and the tool is not cutting into the object, do not include the surface.



968

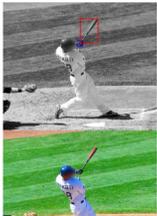
3

For instance, the rolling pin is not touching the counter but is instead only on top of the dough; the rolling pin can't cut through the dough.



Tools that are not touching other objects:

Some of the objects will not be actually touching other objects. In this case, mark nothing to bound. For example, here the baseball bat is in the air. There is nothing to bound. If you think it's possible that the tool is touching another object, please bound the object.



4

Big large objects like floors or lakes or snow:

Sometimes, the tool is touching something like the floor. In this case, please include as much of the floor you can see. Similarly, if a paddle is touching water, include as much of the water as you can see.



969 **H.7 Grasp**

970 Approximately \$3,400 was spent obtaining annotations for grasps over 154K classification tasks. We
971 did a pilot study of thousands of grasp annotations ourselves. This helped us identify a taxonomy that
972 was easily annotated. We then obtained annotations hierarchically.

973 First, grasps were classified into NP-Palm/NP-Fin/Prehensile. Then prehensile grasps were classified
974 into the categories described in Cutkosky [4].

975 **H.8 Prehensile-vs-Non-Prehensile Grasps**

976 In our first round, we obtain annotations of prehensile grasps (where the object is held) compared
977 with two types of non-prehensile grasps: where fingers make contact and where more than the fingers
978 make contact. We frame this as a classification problem.

1

Welcome!

We are trying to classify how humans are interacting with objects with their hands. You will see a double image containing two panels.

In the first, black and white panel, we show a hand in blue and an object in red. Below the image we show the same image in color to help you see better. We are interested in how the hand in blue is interacting with the object in red.



There are three main categories we are interested in are:

1. Prehensile grasps, where the hand is actually holding an object.
2. Non-prehensile fingers only interactions, where the hand is only touching an object (but not holding it) and touching it only with fingers.
3. Non-prehensile more than fingers interactions, where the hand is only touching an object (but not holding it) and touching it with parts of the hand other than the fingers.

Here are several examples of each

3

Non-Prehensile Fingers Only Interactions

These interactions involve the hand touching the object, but not actually holding it.

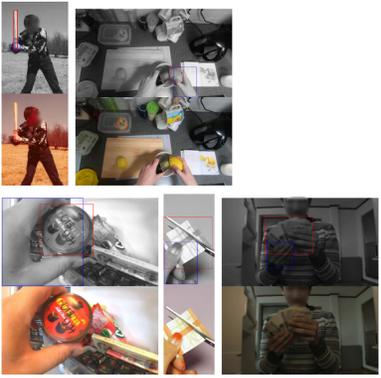
For example:



2

Prehensile grasps

These come in many varieties. What they have in common is that the hand is holding the object and not touching or just supporting the object from below. Holding often requires two directions of force to be applied to the object for instance, one from the fingers, the other from the thumb. A test might be would the hand resist someone grabbing the object? Here are some examples.

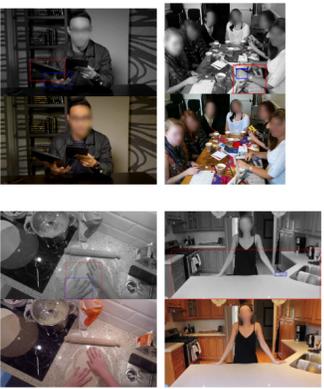


In all cases, the hand is holding the object.

4

Non-Prehensile More than Fingers Interactions

These interactions involve touching objects with parts of the hand that are beyond the fingers. Other parts of the hand include the palm or the backside of the hand. The interaction can involve: fingers and palm, palm only, backside only, fingers and backside. Note that many examples may include supporting from below without holding the object.



980 **H.9 Differentiating Prehensile Grasps**

981 We then obtain annotations for different types of prehensile grasps. The instructions were sub-
982 stantially more complex and the annotation was challenging. We aimed to provide many examples,
983 including illustrations from [12]. Generating fair quality checks was important and we aimed to avoid
984 ambiguous cases.

1

Welcome!

Please note: this task requires careful reading. Guessing will not work.

We are trying to classify how humans are interacting with objects with their hands. You will see a double image containing two panels.

In the first, black and white panel, we show a hand in blue and an object in red. Below the image we show the same image in color to help you see better. We are interested in how the hand is interacting with the object in red.

985



We are looking to categorize human grasps into five categories. There are two key concepts we'll explain that together explain the grasps:

- precision and power grasps
- circular and prismatic grasps

3

986

Power Examples:



5

Prismatic grasps:



987

Prismatic: here the hand is mainly doing a prismatic grasp around a knife (the hand is wrapped around it and the fingers are all in one line). However, the index finger is extended. Please mark this as prismatic.



2

We are looking to categorize human grasps into five categories. There are two key concepts we'll explain that together explain the grasps:

- precision and power grasps
- circular and prismatic grasps

Precision and Power

The first distinction is between precision grasps and power grasps. Precision grasps are used for precise movement, and power is used to exert force on an object. They are distinguished by looking at how much contact happens outside the finger tips and whether the fingers can freely move:

	Amount of contact not on the fingers	Fingers freely moving?
Precision Grasp	The hand is using only the finger tips to contact the object (or at least it's almost entirely the finger tips). The object is not in contact with the palm, and almost all of the contact is at the finger tips.	The fingers can freely move to change where the object is. For instance, a precision grasp of a pen lets you move the pen around easily even if your wrist remains still.
Power Grasp	Parts of the palm of the hand are in contact with the object or most of the contact is not at the finger tips. For instance, if you use all of the fingers to hold an object.	The fingers cannot freely move. Instead, the object is moved around by moving the wrist.

Precision Examples:



4

Circular and Prismatic

The second distinction is between circular grasps and prismatic grasps.

Circular grasps involve finger tips applying pressure on the object from multiple directions. Prismatic grasps involve the tips of fingers 2, 3, 4, and 5 applying pressure from one direction, often in the opposite direction of the thumb. Sometimes three fingers are doing a prismatic grasp and a fourth finger is doing something. In this case, please mark it as a prismatic grasp.

	Shape of the fingers that are in contact	Directions of force of the fingers that are in contact
Circular Grasp	Fingers form a circle around the object. 	The fingers apply force from multiple opposing directions in a circle.
Prismatic Grasp	Fingers form a cylinder or line and apply forces from one side. Note that fingers 2, 3, 4, 5 apply forces in one direction. The thumb opposes this direction. For example: 	The fingers apply force from one direction. Often the thumb applies force from another direction.

Circular grasps:



6

Now we will combine them together

Power Prismatic -- This is a classic power prismatic grasp. The palm is in contact (so it is power). The fingers are wrapped around and are not applying force in multiple opposing directions in a circle.



Here are other ones:



7

Precision prismatic -- here, the fingers are opposed to each other. Just the finger tips are in contact, so it is a precision grasp. Since the in contact fingers are applying force in a single direction and not in a circle around the object, this is a prismatic grasp.



Here are other ones:



988

8

Precision circular -- this is a classic precision circular grasp. The finger tips are wrapped around the object in a circle and are applying forces in multiple directions. But only the finger tips are in contact.



Here are other ones:



9

Power circular -- the palm and lots of the finger are in contact with the pot and the fingers are wrapped around in a circle.



Here are other ones:



989

10

One Exception: Lateral

Most grasps involve using the insides of the fingers. Sometimes people use the side of their fingers. This is an important exception. This is called a lateral grasp where the object is held between the thumb and the side of fingers. This is called a lateral grasp.



Here is another one:



11

Some common issues



Here are some tricky cases and how to resolve them.

Grabbing something non-solid

Sometimes people grab objects that are not solid, such as dough or laundry. Please mark these as power circular (if the hands roughly form a circle or sphere) or power prismatic (if the hands form a cylinder).



990

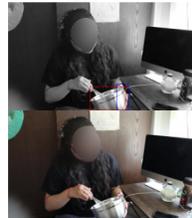
Grabbing bowls

We find these cases difficult. You should label these as power circular (if most of the contact is not from the finger tips) or precision circular (if only the finger tips are in contact).

Grabbing bowls

We find these cases difficult. You should label these as power circular (if most of the contact is not from the finger tips) or precision circular (if only the finger tips are in contact).

For instance, here, the palm is probably in contact with the bowl and certainly a lot of the fingers past the finger tip are too.



13

Funny grasps where fingers are sticking out that don't feel like one category or another

Many grasps do not fall neatly into these categories. Sometimes the hand has an extra finger touching the object. Please use the closest grasp, or the grasp that you think the hand most resembles. For example, this is a power prismatic grasp, but the thumb is in a weird position. We will try to make the best cases as clear examples as possible.



991

992 **H.10 Video Identification**

993 Approximately \$400 was spent annotating thumbnails to provide training data for automatic relevance
 994 filtering. This involved around 21K classification tasks.

995 **H.11 Filtering Videos**

1

996 Welcome!

We are filtering videos from the internet to find real examples where people are interacting with objects. Many of the videos are good, but there are lots of videos that are not suitable for our purposes, and we're asking your help to identify these unsuitable videos so that we can remove them. You will see four thumbnails from a video. These will look like these sorts of images:



We would like you to try to identify three types of videos: not real videos, videos with kids, and videos that are lectures. Everything else is an acceptable video.

3

997 Videos involving children: if any of the thumbnails depict children, or people who appear to be less than 18 years old, we do not want these.

Here are examples of videos involving children. Note that the children may appear in only one of the thumbnails:



5

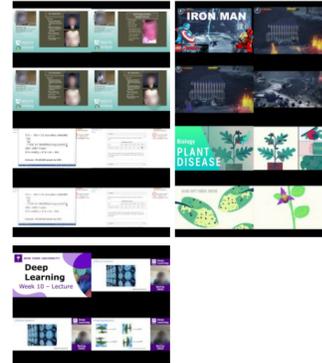
998 If the video doesn't fall into one of these categories, then we're interested in it and you should mark it as Acceptable video. Here are acceptable videos. Note that they are very diverse. This is because we are interested in any video that is not one of the top three categories.



2

Not real videos: these can be cartoons, animations, screen recordings, or slideshows, video games etc. We do not want these.

Here are examples of not real videos. Note that these can be quite diverse since they may be video games, video lectures, cartoons, etc. A video can have a one or two video frames that shows some sort of diagram or logo, but having multiple frames that are not real images makes the video not real.



4

Lecturing/talking videos: some videos depict a person sitting and talking to the camera, and most of the thumbnails have a person (or group of people) in a chair. Here are examples of lecturing videos. If the video has two or more frames where the person is in more or less the same picture, with the same background, and is talking to the camera (in other words, they are not showing something off, like how to disassemble their laptop), please mark it as a lecturing video.



999 **H.12 Counting Hands**

1000 We asked workers to classify videos based on the number of in-contact hands in nine frames from the
 1001 video.

1

1002

Welcome!

We are filtering videos from the internet to find real examples where people are interacting with objects. Many of the videos are good, but there are lots of videos that are not suitable for our purposes, and we're asking your help to identify these unsuitable videos so that we can remove them.

You will see nine frames from a video. These videos are often many minutes long, so we're looking at 9 sample frames from them. These will look like these sorts of images. Each part of the image is frame from the video.

Video 1:

Video 2:

3

1003

There are fewer hands here, but there are at least two (the top-left, middle-left, and middle-middle). This is also a case with at least two frames.

5

1004

Here, you can't see any hands easily. This is also no hands interacting.

1005 **H.13 Identifying Frame Types**

1006 We asked workers to identify up-close vs far-away frames vs non-real images. Gold standard checks
 1007 were chosen here to be deliberately unambiguous.

2

Examples with at least two frames of clearly visible hand interactions.

In these videos two or more frames of the image show hand interaction with an object where the hand is clearly visible. By hand interaction we mean that a person is picking up, holding, or touching some object with the intent of doing something with that object. If the hand is in contact with an object (basically anything other than a person), then it counts as interaction with the one exception of resting hands on a table or desk or chair. By clearly visible we just mean that you do not have to look carefully to see if someone in the background is maybe holding something.

Here are some examples:

Here, there are lots of hands interacting with things.

Here there are also lots of hands interacting with things

4

Examples with fewer than two frames of visible hand interactions (including no frames of hand interaction).

These will mainly be examples where you can't see hands or the people are sitting and talking to the camera. They are often pretty easy to see.

Here, two people are talking in this video. You can see a hand if you look really carefully, but it's not holding anything. This clearly does not show a lot of interaction.

Here, there is a hand visible in one frame. This also shows fewer than two hand interactions.

1

Welcome!

We're looking to categorize frames in YouTube videos according to how close to the camera the objects are. We have three main categories based on where most of the objects in the scene are. If there are hands visible, please use the hands that you can see to decide the distance.

1. Within 0.5m of the camera.

These are frames where the hands or other objects are very close to the camera and take up much of the frame. These image should feel as if the objects could hit you.



1008

2

2. More than 1m away

These are images where a lot of the scene is visible, and if there are hands, they're probably 1m away. If the object is between 50cm and 1m, please make your best guess as to which is closer. The qualifier and other tests will only have ones where the distance is obvious.



3. Not a real images

These are images that are not actually images, since they are things like text or diagrams. If it doesn't look like a real image, please mark it as "not a real image".



1009 **H.14 Image Redaction**

1010 Approximately \$4800 was spent redacting children and faces from the dataset, spent over approxi-
1011 mately 500K tasks that were primarily binary classification, but also included approximately 17K
1012 box annotation tasks.

1013 In this section, we indicate unblurred faces with a black and white checkerboard pattern, and hide
1014 minors with a black mask.

1015 **H.14.1 Instructions for Unblurred Face Spotting**

1

Welcome!
We're looking to spot cases where an automatic face blurring system failed. We are classifying images into two categories:

- All faces blurred: all the visible faces that are in the image are blurred
- Some faces unblurred: there is at least one image in which there is an unblurred faces

All Faces Blurred
Here are some examples of the images where all the faces are blurred



1016

2

Note that you cannot see the person's face, so there is no visible unblurred faces.



3

Some Faces Unblurred
Here are some examples where some of the faces are not blurred



1017



4



1018 **H.14.2 Instructions for Unblurred Face Bounding**

1

Welcome!
 We're trying to spot faces that we missed while blurring all the faces
 Instructions:
 • Please draw a bounding box around any face that is visible and is not blurred.
 • The box should go around any visible parts of the face. The part of the face that we are interested in covering is shown in green below and should definitely include the ears, nose, eyes, mouth, and chin and forehead. It does not need to include the neck.

1019



2



3

Common Case 2:
 Many times the faces will be in the background. You only need to draw boxes around faces that are visible. If the person is not facing the camera and no face is visible, then you do not have to draw a box around that face. For instance, in this image there is one person who has a face that is unblurred and two people whose faces are not visible. Please draw a box only around the person whose face is visible and not blurred.

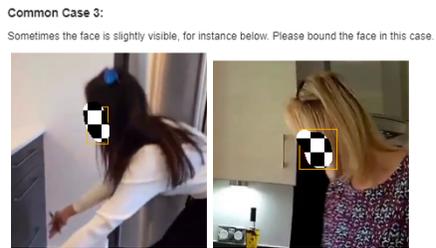
1020



If the face is too small, for instance a person that is only 10 pixels high total, you do not need to mark the face. For instance, this image, there are two people far away on the right. They are far too far away to be recognized and so their faces are already blurred. If a person is this size or smaller, you can skip them.



4



1021 **H.14.3 Instructions for Spotting Minors**

1022 Workers were instructed to spot children in data and were told that in unsure cases (e.g., someone in
 1023 their teens), they should mark that person as a child.

1

2



Welcome!

We are looking to see if there are people who appear to be younger than 18 years of age. The two categories are:

- Shows children: It appears that there is at least one child in the image
- No children: It appears that there are no children in the image

1024

We are defining children as people who are under 18. In some cases, this may be difficult to tell. If you are not sure, but think it's fairly likely that the person is a child, please mark them as a child. For example, if you're positive the person is 17-18, but not sure where, please mark them as a child.

Examples of the children category

Here are examples where there are people who appear to be younger than 18 in the image.



Please note that sometimes children may not have their faces visible. Please label these as children. For example, in this image, the person is clearly a child.



3

4

Examples of the no children category

Here are examples of no children category



1025



If you cannot recognize whether the person is a child (e.g., a face in a distant crowd that you can't actually see), you can assume the person is an adult. If you can't quite see the person but it seems likely the person is a child (e.g., based on clothing or size), you should mark the person as a child. We will only test you on examples that we believe are clearcut.

1026 **H.15 Polygon Labeling**

1027 Approximately \$400 was obtaining polygons for comparison with SAM for both hands and objects.
1028 This was done across 2000 tasks. The two tasks are explained below, and follow as similar pattern to
1029 other annotations done in the dataset: a top images illustrates what is to be annotated, and a bottom
1030 image is annotated.

1031 **H.15.1 Hands**

1

Welcome!



We are asking you to segment a hand in the image. Do not segment any object other than a hand.
If you've tried this before, we have fixed the bug that prevented passing the qualifier.
We will identify the hand that you need to segment by providing a box around it in the image.
• Please draw a segment that entirely covers the hand in the image.
• The segment should never go outside the box.
• The box is a light box around the object, so the segment will almost always be inside the box at some point.
• Please segment the hand in the bottom half of the image. The top half of the image is only to help you identify the hand.
For instance, the hand is identified here in the top box in the image. The segment of the hand is drawn below.

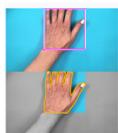


2

Easy Cases



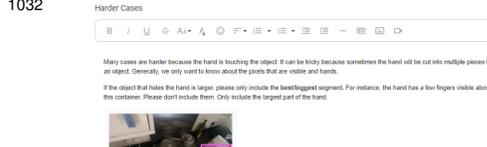
Some cases are easy and clear because the hand is not contacting an object.
For these cases, the only difficulty is figuring out where the hand ends. Please use either the wrist or the extent of the hand in the image.
Notice that the image changes from colored to black and white at the edge of the box. This may help you identify where to stop the segment.



1032

3

Harder Cases



Many cases are harder because the hand is touching the object. It can be tricky because sometimes the hand will be cut into multiple pieces by an object. Generally, we only want to know about the parts that are visible and hands.
If the object that hides the hand is larger, please only include the best/biggest segment. For instance, the hand has a few fingers visible above the container. Please don't include them. Only include the largest part of the hand.



4

However, sometimes an object is in front of the hand and is small but cuts the hand into multiple pieces. If the object splitting the hand into multiple pieces is small (about the thickness of a pen), please pretend the object isn't there.



However, note that while the segment does include the card, it does not include the object the hand is holding.

5

Two More Examples



Here are two more examples.

References

- 1036 [1] The vision for intelligent vehicles and applications (VIVA) challenge, laboratory for intelligent and safe
1037 automobiles, UCSD. <http://cvrr.ucsd.edu/vivachallenge/>.
- 1038 [2] Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing
1039 mutual information across views. *arXiv preprint arXiv:1906.00910*, 2019.
- 1040 [3] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial
1041 gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR,
1042 2018.
- 1043 [4] Mark R Cutkosky et al. On grasp choice, grasp models, and the design of hands for manufacturing tasks.
1044 *IEEE Transactions on robotics and automation*, 5(3):269–279, 1989.
- 1045 [5] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, , Antonino Furnari, Jian Ma, Evangelos Kazakos,
1046 Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Rescaling egocentric
1047 vision. *CoRR*, abs/2006.13256, 2020.
- 1048 [6] Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar, Richard Higgins, Sanja Fidler, David
1049 Fouhey, and Dima Damen. Epic-kitchens visor benchmark: Video segmentations and object relations. In
1050 *Proceedings of the Neural Information Processing Systems (NeurIPS) Track on Datasets and Benchmarks*,
1051 2022.
- 1052 [7] David F. Fouhey, Weicheng Kuo, Alexei A. Efros, and Jitendra Malik. From lifestyle vlogs to everyday
1053 interactions. In *CVPR*, 2018.
- 1054 [8] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE
1055 international conference on computer vision*, pages 2961–2969, 2017.
- 1056 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
1057 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 1058 [10] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
1059 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint
1060 arXiv:2304.02643*, 2023.
- 1061 [11] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
1062 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on
1063 computer vision*, pages 740–755. Springer, 2014.
- 1064 [12] Jia Liu, Fangxiaoyu Feng, Yuzuko C Nakamura, and Nancy S Pollard. A taxonomy of everyday grasps in
1065 action. In *2014 IEEE-RAS International Conference on Humanoid Robots*, pages 573–580. IEEE, 2014.
- 1066 [13] A. Mittal, A. Zisserman, and P. H. S. Torr. Hand detection using multiple proposals. In *BMVC*, 2011.
- 1067 [14] Supreeth Narasimhaswamy, Thanh Nguyen, Mingzhen Huang, and Minh Hoai. Whose hands are these?
1068 hand detection and hand-body association in the wild. In *Proceedings of the IEEE/CVF Conference on
1069 Computer Vision and Pattern Recognition (CVPR)*, 2022.
- 1070 [15] Supreeth Narasimhaswamy, Zhengwei Wei, Yang Wang, Justin Zhang, and Minh Hoai. Contextual attention
1071 for hand detection in the wild. In *Proceedings of the IEEE/CVF international conference on computer
1072 vision*, pages 9567–9576, 2019.
- 1073 [16] Shengyi Qian, Linyi Jin, Chris Rockwell, Siyi Chen, and David F. Fouhey. Understanding 3d object
1074 articulation in internet videos. In *CVPR*, 2022.
- 1075 [17] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,
1076 Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge.
1077 *International journal of computer vision*, 115(3):211–252, 2015.
- 1078 [18] Dandan Shan, Jiaqi Geng, Michelle Shu, and David F Fouhey. Understanding human hands in contact at
1079 internet scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
1080 pages 9869–9878, 2020.
- 1081 [19] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu,
1082 Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition.
1083 *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3349–3364, 2020.
- 1084 [20] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019.
- 1085 [21] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transforma-
1086 tions for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern
1087 recognition*, pages 1492–1500, 2017.
- 1088 [22] Kaiyu Yang, Jacqueline Yau, Li Fei-Fei, Jia Deng, and Olga Russakovsky. A study of face obfuscation in
1089 ImageNet. In *International Conference on Machine Learning (ICML)*, 2022.
- 1090