
Supplementary Material for IKEA Manuals at Work: 4D Grounding of Assembly Instructions on Internet Videos

Yunong Liu¹ Weiyu Liu¹ Shubh Khanna¹ Cristobal Eyzaguirre¹ Manling Li¹
Juan Carlos Niebles¹ Vineeth Ravi² Saumitra Mishra² Jiajun Wu¹

¹Stanford University ²J.P. Morgan AI Research
Website: <https://yunongliu1.github.io/ikea-video-manual/>

1 Contents

2	A Dataset Details	1
3	B Datasheets	3
4	C Details for Data Annotation	5
5	D Details for Quality Control	8
6	E Annotation Interfaces and Instructions	9
7	F Error Bar	11
8	G Compute Resources	11

9 A Dataset Details

10 IKEA Video Manuals is a large-scale multimodal dataset with high-quality, spatial-temporal align-
11 ments of step-by-step instructions, 3D object representations, and real-world video demonstrations
12 from the Internet. IKEA Video Manuals provides 34,441 annotated video frames, aligning 36 IKEA
13 manuals with 98 assembly videos for six furniture categories. Fig. A1 shows all 3D furniture models
14 included in the dataset. An example of the annotations associated with each frame is shown in Fig. A2.
15 We provide details of the data and annotations associated with each frame below.

16 Furniture-level information

- 17 • **Category:** The category label of the furniture (e.g., Bench).
- 18 • **Name:** The furniture name (e.g., applaro).
- 19 • **Furniture IDs:** A list of IKEA product IDs for the furniture.
- 20 • **Variants:** A list of furniture variants, if applicable.
- 21 • **Furniture URLs:** A list of IKEA product page URLs for the furniture.
- 22 • **Furniture Main Image URLs:** A list of URLs for the main product images on the IKEA website.

23 Video-level information

- 24 • **Video URL:** The URL of the video.
- 25 • **Additional Video URLs:** A list of additional video URLs for the same furniture.
- 26 • **Title:** The title of the video.
- 27 • **Duration:** The duration of the video (in seconds).



Figure A1: All furniture items included in the IKEA Video Manuals dataset, categorized by type—Desk, Table, Chair, Bench, and Misc.

- 28 • **Resolution:** The resolution of the video (e.g., 1920x1080).
- 29 • **FPS:** The frame rate of the video (e.g., 30).
- 30 • **People Count:** The number of people in the video.
- 31 • **Person View:** The view of the person in the video (e.g., front, side).
- 32 • **Camera Fixation:** The fixation of the camera in the video (e.g., static, moving).
- 33 • **Indoor/Outdoor Setting:** The setting of the video (e.g., indoor, outdoor).

34 **Assembly step information**

- 35 • **Step ID:** An unique ID is assigned to the assembly step.
- 36 • **Step Start:** The start time of the assembly step is shown in the video.
- 37 • **Step End:** The end time of the assembly step is shown in the video.
- 38 • **Substep ID:** The unique ID assigned to the substep within the assembly step.
- 39 • **Substep Start:** The start time of the substep in the video.
- 40 • **Substep End:** The end time of the substep in the video.

41 **Frame-level information**

- 42 • **Frame Time:** The timestamp of the frame in the video.
- 43 • **Number of Camera Changes:** The number of camera changes that have been labelled before the
- 44 current frame.
- 45 • **Frame Parts:** A list of parts that are labelled in the frame (e.g., $[\{0, 2\}, \{1\}, \{3\}]$). The sub-
- 46 assemblies that have been constructed in previous steps are denoted by a tuple of the part IDs.
- 47 • **Frame ID:** An unique identifier for the frame (e.g., 1584).
- 48 • **Is Keyframe:** A boolean value indicating whether the frame is a keyframe.
- 49 • **Is Frame Before Keyframe:** A boolean value indicating whether the frame is immediately before
- 50 a keyframe.
- 51 • **Frame Image:** The RGB image of the frame.
- 52 • **Annotated Masks:** A list of segmentation masks for the parts in the frame.
- 53 • **Annotated Poses:** A list of 6D poses for the parts in the frame.

54 **Manual information**

- 55 • **Manual Step ID:** A unique ID assigned to the assembly step. This ID is associated with the step
- 56 IDs in frame annotations.
- 57 • **Manual URLs:** A list of URLs for the assembly manual PDFs.
- 58 • **Manual ID:** An unique ID of the assembly manual from IKEA.
- 59 • **Manual Parts:** A list of part IDs is shown in the corresponding manual step.
- 60 • **Manual Connections:** List of connections between parts in the manual step.
- 61 • **PDF Page:** The page number of the manual step in the PDF.
- 62 • **Cropped Manual Image:** The cropped image of the corresponding manual step.

63 **B Datasheets**

64 **Dataset description.** The datasheet for the IKEA Video Manuals dataset, available at <https://github.com/yunongLiu1/IKEA-Manuals-at-Work/blob/main/datasheet.md>, including

65 key aspects of data collection and annotation:

66

- 67 • **Consent:** The dataset is built upon two existing datasets, IKEA-Manual and IAW, which are
- 68 publicly available for research purposes under the Creative Commons Attribution 4.0 International
- 69 (CC-BY-4.0) license.
- 70 • **Personally Identifiable Information and Offensive Content:** The dataset does not contain any
- 71 personally identifiable information or offensive content, as it focuses on furniture objects and
- 72 assembly instructions.

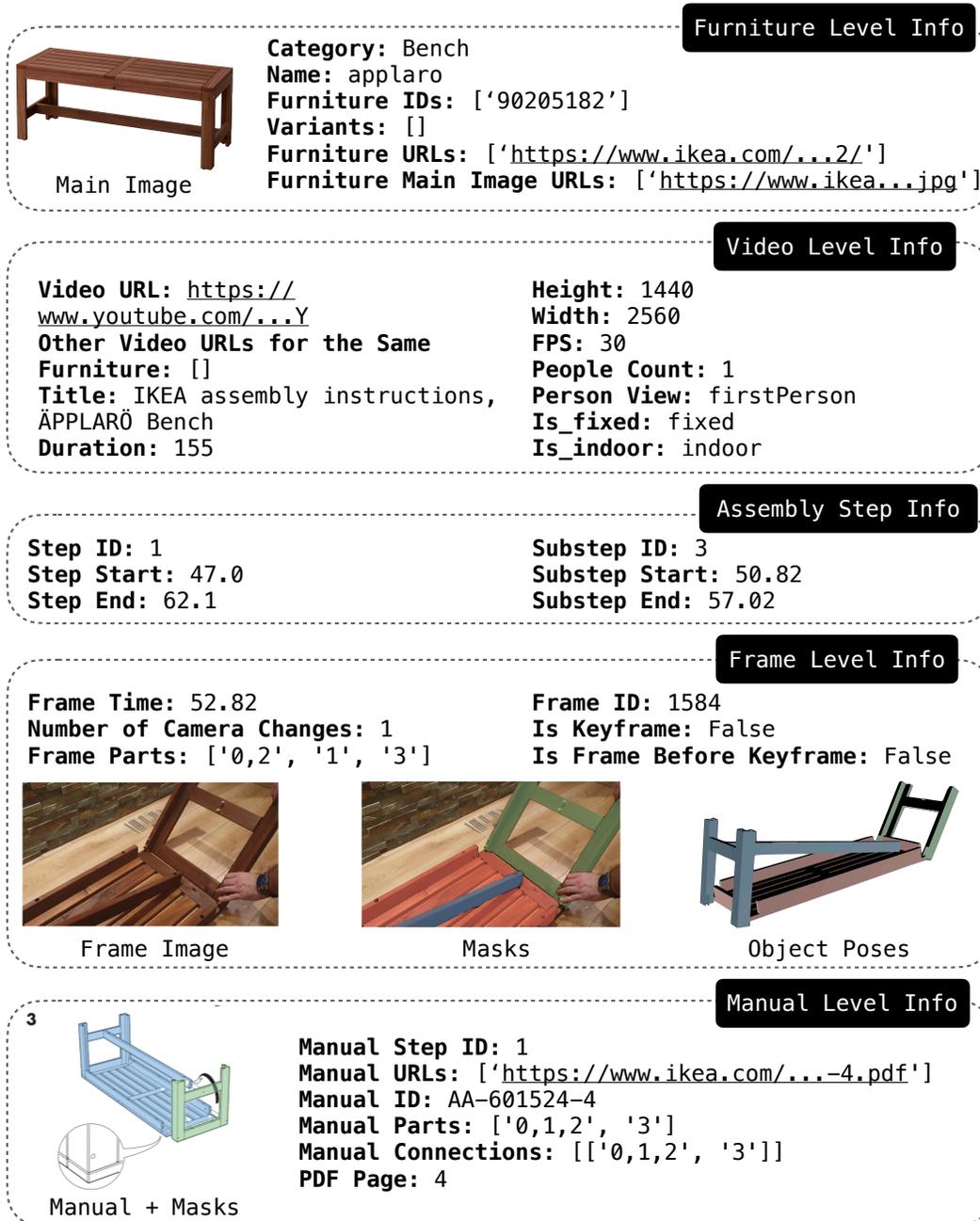


Figure A2: An example of an annotated frame in the IKEA Video Manuals dataset. The annotation and data are divided into furniture-level, video-level, assembly step-level, frame-level, and manual-level information.

- 73 • **Annotation Process and Compensation:** The data annotation process was outsourced to an
- 74 annotation company. The annotators were compensated based on the work they provided, with the
- 75 estimated hourly pay being above the minimum wage.
- 76 Please refer to the datasheet for more detailed information on the dataset.
- 77 **Link and license.** The dataset is available for public access under the CC-BY-4.0 license: <https://github.com/yunongLiu1/IKEA-Manuals-at-Work>
- 78

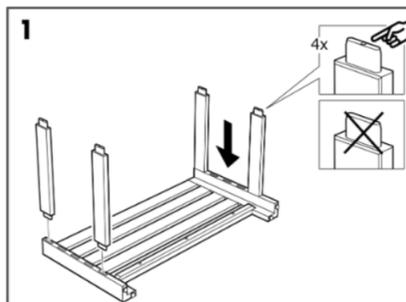


Figure A3: An example of an assembly step from the IKEA instruction manual that involves the assembly of four parts.

79 **Maintenance.** The dataset is hosted on GitHub and will be maintained by the authors. The repository
 80 can be found at: <https://github.com/yunongLiu1/IKEA-Manuals-at-Work>. The dataset has
 81 the following DOI: <https://doi.org/10.5281/zenodo.11623997>

82 **Author statement.** The authors bear all responsibility in case of violation of rights. All annotations
 83 were collected by the authors and we are releasing the dataset under CC-BY-4.0.

84 **Format.** The dataset contains videos, 3D models, manual PDFs, and annotated data (including tem-
 85 poral step alignments, temporal substep alignments, 2D-3D part correspondences, part segmentations,
 86 part 6D poses, and estimated camera parameters). The annotated data is stored in the JSON format.
 87 Other data are stored in their original formats, and uploaded in a zip file. Upon decompression, the
 88 dataset is organized into subdirectories for videos, 3D models, and manual PDFs. Each is organized
 89 into subdirectories for furniture categories, and further subdirectories for individual furniture items.

90 **Croissant Metadata** We will provide the structured metadata (schema.org standards) in <https://github.com/yunongLiu1/IKEA-Manuals-at-Work/metadata.json>.

92 C Details for Data Annotation

93 To create the IKEA Video Manuals dataset, we identified 36 IKEA objects from the IKEA-Manual
 94 dataset [1] that have corresponding assembly videos in the IAW dataset [2]. We matched the unique
 95 IDs of the instruction manuals to ensure correct correspondence between the datasets. We provide
 96 additional details for each of the annotation steps below.

97 C.1 Annotating Assembly Steps

98 For each assembly step annotated in the IKEA-Manual dataset [1], we identify matching video
 99 segments in the IAW dataset [2]. We manually adjust the start and end time of each video segment to
 100 include a more complete assembly process, from picking up a part to positioning and tightening. The
 101 adjustment ensures better alignment with the physical assembly actions.

102 C.2 Annotating Assembly Substeps

103 In the IKEA instruction manuals, each single step may involve the assembly of multiple parts (as
 104 shown in Fig. A3). We provide a more fine-grained assembly process by introducing *substeps*. A
 105 substep is labelled when 1) a new part appears in the video or 2) a new sub-assembly is created
 106 through positioning and/or fastening of parts. On average, each assembly step contains 7.59 substeps.
 107 In total, the IKEA Video Manuals dataset contains 1120 substeps.

108 C.3 Annotating Part Identities

109 In our dataset, each part of the 3D furniture model is assigned a unique ID consistent with the
 110 IKEA-Manual dataset. However, locating individual 3D furniture part in the frame can be challenging
 111 due to several ambiguities, as illustrated in Fig. A4:

- 112 (a) Wrongly assembled parts that are initially placed incorrectly and later relocated, causing
 113 confusion for the annotator (Fig. A4a).

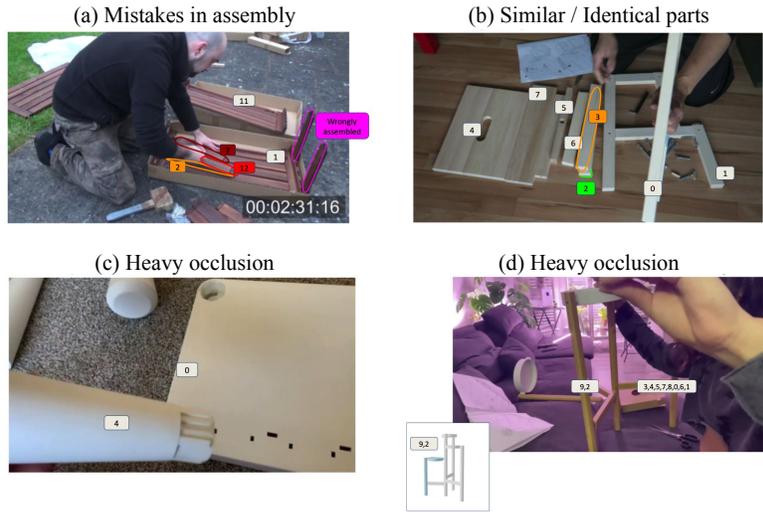


Figure A4: Examples of ambiguities in annotating part identities in assembly videos: (a) Wrongly assembled parts later relocated, (b) Similar-looking parts that are difficult to distinguish, (c-d) Heavily occluded parts and boundaries between parts.



Figure A5: An example of refining the part segmentation using the brush tool. The initial segmentation is generated by the Segment Anything Model (SAM) model.

114 (b) Similar or identical-looking parts, such as chair legs, that are difficult to distinguish and
 115 label accurately (Fig. A4b).

116 (c) Parts that are heavily occluded, making it challenging to recognize the parts and their
 117 boundaries (Fig. A4c-d).

118 To ensure accurate part tracking throughout the video, we manually label the parts in the first frame
 119 of each substep after watching the entire video. This annotation is crucial for maintaining consistency
 120 when annotators only see individual frames in subsequent annotations. This approach ensures
 121 consistent part identities throughout the video, addressing challenges posed by heavy occlusions,
 122 similar-looking parts, and assembly mistakes.

123 C.4 Annotating Segmentation Mask

124 To efficiently generate segmentation masks for furniture parts in video frames, we utilize the Segment
 125 Anything Model (SAM). When SAM fails to generate accurate masks (e.g., Fig. A5), we use a brush
 126 tool built into our annotation interface to refine the masks manually.

127 C.5 Annotating 2D-3D Keypoints

128 To establish correspondence between the 3D parts and their 2D projections in the video frames,
 129 we annotate keypoints on both the 3D parts and the 2D images. Our annotation interface is shown

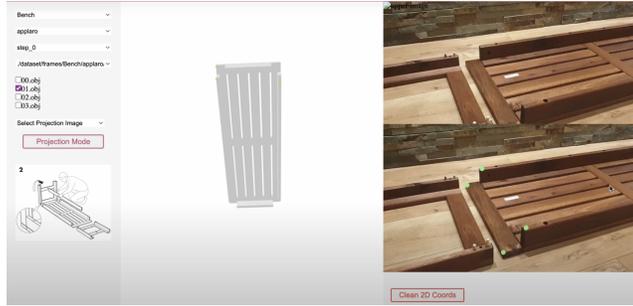


Figure A6: The annotation interface for labelling keypoint correspondences between 3D models and 2D video frames.

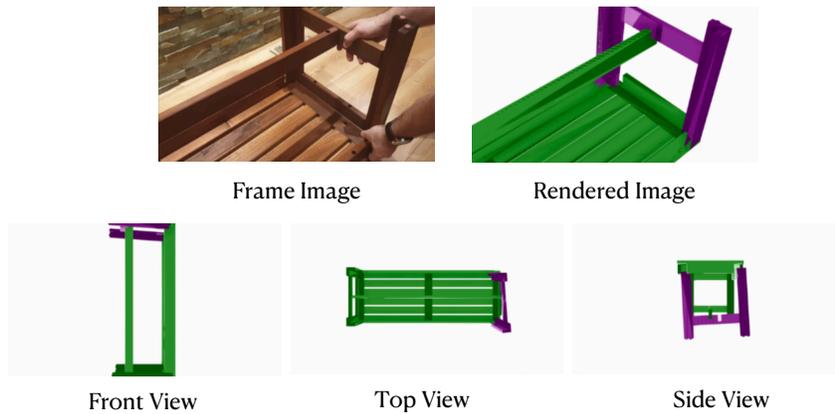


Figure A7: An example of an inaccurate part pose estimated from the 2D-3D keypoints. Despite that the 2D projections of the 3D models overlap with the parts in the video frame (top row), the 3D poses of the parts can be found incorrect when visualized in 3D and viewed from other angles (bottom row).

130 in Fig. A6. The annotation interface computes the part poses and camera parameters using the
 131 Perspective-n-Point (PnP) algorithm and visualizes the 2D projection in real-time. Based on the
 132 visualization, the annotator can interactively refine the keypoint annotations to maximize the overlap
 133 between the 2D projection and the part seen in the 2D image.

134 C.6 Annotating Camera Changes

135 A prerequisite for achieving spatially and temporally accurate pose annotation is a correct estimation
 136 of camera parameters from the video. Many videos in the dataset include changes in camera
 137 viewpoints, camera movements, and adjustments of focal lengths. These changes can potentially lead
 138 to different camera parameters. To account for these factors, we manually annotate camera changes
 139 in the IKEA Video Manuals dataset. By annotating all frames when a camera change occurs, we can
 140 estimate the intrinsic parameters for each video segment between two camera changes, assuming that
 141 the intrinsic parameters remain consistent within the segment.

142 C.7 Pose Refinement

143 While initial estimates of the part poses can be obtained from the annotated 2D and 3D keypoints,
 144 these estimates are often inaccurate, particularly in terms of the relative positions and orientations of
 145 the parts in 3D space. Fig. A7 shows an example where the 2D projection of a part appears correct,
 146 but when viewed in 3D, the part is positioned incorrectly relative to other parts. To address this
 147 issue, we developed an interface that allows annotators to refine the initial estimate by rotating and
 148 translating each part in 3D space. The annotators can view 3D parts from different viewpoints to
 149 ensure that the relative poses of the parts are correct and consistent with the assembly process seen in
 150 the videos.

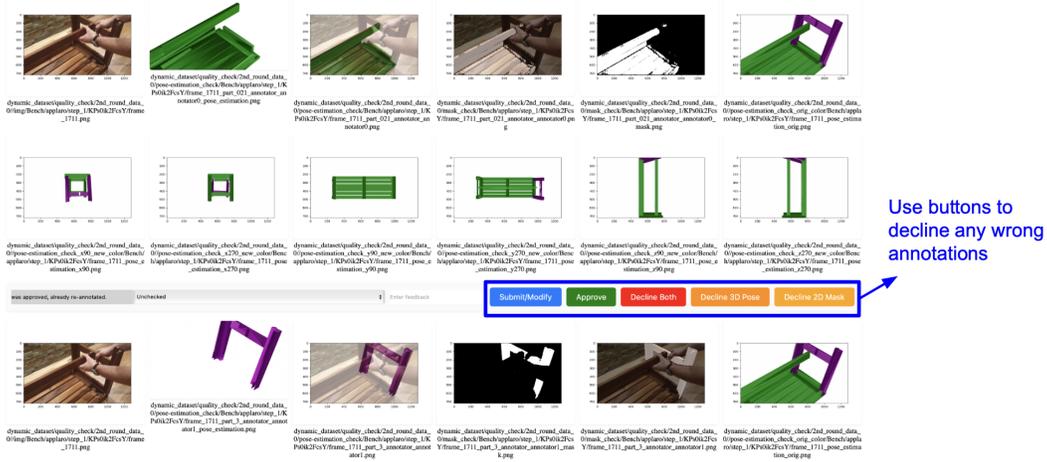


Figure A8: The verification interface for assessing the quality of mask and pose annotations. The interface visualizes the original video frame, mask overlay, pose overlay, and 3D parts in estimated poses from different viewpoints.

151 We take an additional step in the pose refinement process to help maintain the temporal smoothness
 152 of the part trajectories. During the pose refinement process, the initial poses of the parts in a frame are
 153 set to the refined poses of the corresponding parts from the previous frame. This initialization strategy
 154 helps to reduce large changes in the annotated part poses between neighboring frames, resulting in
 155 more coherent and realistic pose trajectories.

156 D Details for Quality Control

157 To ensure the accuracy and consistency of the annotations in the IKEA Video Manuals dataset, we
 158 analyze the common errors in the annotations and perform extensive verifications.

159 D.1 Common Errors

160 For mask annotations, common errors include incorrect part segmentation, missing parts, and noisy
 161 masks. Incorrect part segmentation occurs when annotators misidentify the boundaries of a part due
 162 to similar colours or complex shapes. Missing parts occur when certain parts are not segmented,
 163 often due to occlusions. Noisy masks often occur when the SAM model fails to generate accurate
 164 masks, leading to incomplete or inaccurate segmentation.

165 For pose annotations, common errors include incorrect part identification, incorrect relative poses,
 166 and interpenetrations. Incorrect part identification occurs when the annotators annotate an incorrect
 167 part, leading to an incorrect pose. Incorrect relative poses occur when the estimated pose does not
 168 accurately reflect the actual position and orientation of the part relative to other parts in 3D space.
 169 Interpenetrations occur when parts intersect or overlap in 3D space, leading to unrealistic poses.

170 D.2 Extensive Verification

171 We conduct extensive verification to ensure the high quality of the mask and pose annotations. The
 172 verification interface (as shown in Fig. A8) displays the original video frame, the video frame overlaid
 173 with segmentation masks, the video frame overlaid with 2D projections based on estimated poses, and
 174 the 3D parts in the estimated poses from different viewpoints. In particular, for mask annotations, we
 175 verify if the 2D mask corresponds to the correct part, covers the entire part, does not contain pixels of
 176 other parts, and is free of noise due to limitations of the Segment Anything Model (SAM). For pose
 177 annotations, we verify if the pose annotation corresponds to the correct part, the 2D projection aligns
 178 with the part in the frame image, and the 3D parts have correct relative poses. We automatically filter
 179 out pose annotations with interpenetrations between parts.

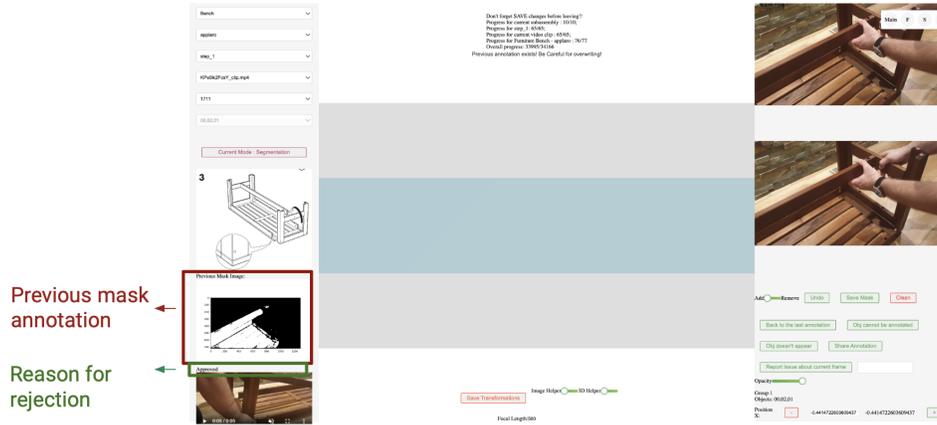


Figure A9: The interface for correcting and refining mask annotations based on feedback. The interface provides tools for manual refinement of the segmentation masks including a brush and an eraser.

180 E Annotation Interfaces and Instructions

181 This section provides details of the instructions given to annotators. Instructions for using these inter-
 182 faces are mainly provided through demonstration videos, which are included in the project's GitHub
 183 repository (<https://github.com/yunongLiu1/IKEA-Manuals-at-Work>) for reference.

184 E.1 Segmentation Mask and 2D-3D Points Correspondence Annotation Interface (Fig. A6)

185 The Segmentation Mask and 2D-3D Points Correspondence Annotation Interface allow annotators
 186 to generate segmentation masks and establish correspondences between 3D models and 2D video
 187 frames. Annotators can switch between the two annotation modes using a dedicated button in the
 188 interface. The following steps outline the annotation process:

- 189 1. Select the appropriate category, subcategory, object, and step for the video you want to
 190 annotate.
- 191 2. In the Segmentation Mask mode:
 - 192 • Select points that best represent the overall shape and area of the part to ensure optimal
 193 performance of the Segment Anything Model (SAM).
 - 194 • Use the provided tools, such as a brush or eraser, to refine the mask based on the
 195 feedback provided.
- 196 3. In the 2D-3D Points Correspondence mode:
 - 197 • Select corresponding points on the 3D model and the 2D video frame that represent
 198 key features or edges of the furniture parts.
 - 199 • Review the rendered image and adjust the selected points if necessary to improve the
 200 alignment between the 3D model and the 2D video frame.
- 201 4. Navigate frames using the 'Next Frame' button and review the predicted points from the
 202 TAPIR model, modifying any unsatisfactory points.
- 203 5. Review the segmented images and estimated poses for accuracy and consistency, and submit
 204 the annotations.

205 E.2 Mask Re-Annotation Interface (Fig. A9)

- 206 1. **Review Previous Mask:**
 - 207 • The interface will display the previously annotated mask in the bottom left corner of
 208 the screen, along with the reason for the decline. Reviewing the previous mask and the
 209 reason for the decline helps the annotator understand the required corrections.

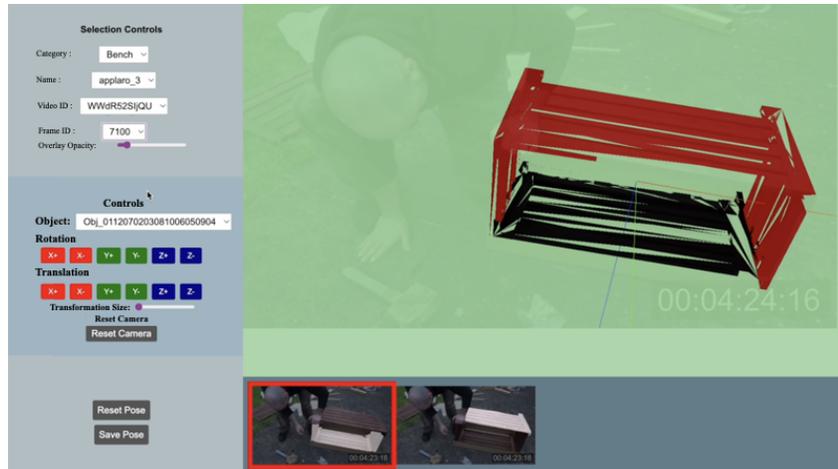


Figure A10: The interface for manually adjusting part poses in 3D. The interface supports adjusting the 3D position and orientation of each part. The interface also provides visualization of the 3D parts from different perspectives.

- 210
- 211
- 212
- 213
- Reasons for the decline include "mask was annotated to the wrong part," "mask did not include the whole part/include other parts," or "noisy mask, which is normally caused by the limitation of SAM and can be solved by using the brush." These specific reasons guide the annotator in refining the mask.

214 **2. Refine Mask:**

- 215
- 216
- 217
- 218
- 219
- 220
- 221
- Use the provided tools, such as a brush or eraser, to refine the mask based on the feedback provided. These tools allow precise modifications to the mask.
 - Ensure the refined mask accurately captures the entire part while excluding any neighbouring parts or background. An accurate and complete mask is essential for downstream tasks.
 - Pay attention to the edges and boundaries of the part to create a clean and precise mask. Well-defined edges improve the quality and usability of the mask.

222 **3. Additional Buttons:** (Same as in the Segmentation Mask Annotation Interface)

- 223
- 224
- 225
- 226
- 4. **Review and Submit:** Review the refined mask for accuracy and completeness, ensuring it addresses the reason for the decline. This final review step verifies that the necessary corrections have been made. Submit the updated mask using the provided submission functionality to save the work and proceed to the next task.

227 **E.3 Pose Refinement Interface (Figure A10)**

228 The Pose Refinement Interface enables annotators to refine the initial poses estimated from the
229 previous annotation. The following steps outline the pose refinement process:

- 230
- 231
- 232
- 233
- 234
- 235
1. Review the initial poses of all parts in the frame, estimated from the previous annotation.
 2. Use the provided controls to adjust the position and orientation of each part in the camera frame.
 3. Ensure that the relative positions and orientations of the parts are consistent with the assembly process.
 4. Review the refined poses for accuracy and submit the updated poses.

236 By following these instructions and leveraging the provided video demonstrations, annotators can
237 effectively use the annotation interfaces to generate high-quality segmentation masks, 2D-3D point
238 correspondences, and refined part poses for the IKEA Video Manuals dataset.

239 **F Error Bar**

240 To assess the variability of our model’s performance, we run experiments on a subset of the data with
241 3 different random seeds for both the segmentation and pose estimation tasks. For part-conditioned
242 segmentation, the standard deviation of the IoU metric is 0.01 for the CNOS method and 0.03 for
243 the SAM-6D method. When considering the Top-5 IoU, the standard deviations are 0.08 and 0.09
244 respectively. For part-conditioned 6D pose estimation, the SAM-6D method has a standard deviation
245 of 0.12 for the ADD score and 0.08 for the ADD-S score. The MegaPose method has standard
246 deviations of 0.09 and 0.05 for ADD and ADD-S. These results indicate that the performance of
247 these models on our dataset remains relatively consistent overall.

248 **G Compute Resources**

249 The computational resources used for this project were computing nodes from the Stanford SC
250 computational cluster. We used around 40 jobs lasting 7 days for running segmentation experiments
251 using SAM-6D and CNOS, and pose estimation experiments using SAM-6D and MegaPose. The
252 jobs were assigned to nodes equipped with different NVIDIA GPU models, including 2080 Ti, Titan
253 RTX, 3090, A40, A5000, A6000, and L40S, based on availability.

254 **References**

- 255 [1] Wang, R., Zhang, Y., Mao, J., Zhang, R., Cheng, C.Y., Wu, J.: Ikea-manual: Seeing shape
256 assembly step by step. *Advances in Neural Information Processing Systems* **35**, 28428–28440
257 (2022)
- 258 [2] Zhang, J., Cherian, A., Liu, Y., Ben-Shabat, Y., Rodriguez, C., Gould, S.: Aligning step-by-step
259 instructional diagrams to video demonstrations. In: *Conference on Computer Vision and Pattern*
260 *Recognition (CVPR)* (2023)