

---

# AFUN: Towards an Affordance Foundation Model for Functionality Understanding

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Affordance understanding bridges visual perception and physical action, serving  
2 as an explainable interface for robot manipulation in open and unstructured real-  
3 world environments. Yet, building an affordance foundation model that not only  
4 understands *where* and *how* the interaction should happen, but also generalizes  
5 across diverse environments, objects, and tasks, remains a long-standing research  
6 challenge. Existing methods typically address only part of this challenge, either  
7 localizing task-relevant regions without specifying executable motion, or predicting  
8 motion but with limited scalability. In this paper, we present **AFUN**, a step  
9 towards an affordance foundation model for functionality understanding. From a  
10 single RGB-D observation and a language task description, AFUN predicts a task-  
11 conditional functional mask (*where* to interact) and a 3D post-contact motion curve  
12 (*how* to interact). To support open-world generalization, we build a large-scale  
13 standardized data pipeline that converts heterogeneous robot, human, simulation,  
14 and real-world scan data into a shared affordance schema with language, masks,  
15 and object-centric 3D motion labels. We evaluate AFUN from three aspects: for  
16 affordance segmentation, AFUN outperforms all baselines by a large margin across  
17 8 test sets from 4 benchmarks, improving mean gIoU/cIoU by **+23.9/+26.3**; for  
18 contact-point prediction, it predicts substantially more accurate points, with a  
19 **12.7–61.3%** hit-rate gain over the best baseline; and for 3D motion, it achieves the  
20 best performance on all three test sets. AFUN can be deployed for real-world robot  
21 manipulation without finetuning for robot embodiment, demonstrating the ability  
22 to adapt to open-world affordance tasks.

## 23 1 Introduction

24 Imagine stepping into a brand-new bedroom; a human can already understand *which* object can  
25 do *what*, and *how* to do it. For instance, a drawer can be opened or closed from its handle, and  
26 a human can identify the exact grasping location before the actual action. This concept of visual  
27 *affordance* [20] to understand objects’ functionalities underpins human’s capability to perform daily  
28 tasks in unstructured real-world environments [69, 62]. In robotics and embodied AI, affordance  
29 understanding serves as a crucial and explainable interface between visual understanding and physical  
30 action. Yet, building a foundation model for affordance understanding that can scale across diverse  
31 environments, objects, and tasks is a long-standing research challenge.

32 There are three interconnected requirements when building such an affordance foundation model.  
33 **(I)** The dataset used to train the model must reflect the diversity of real-world manipulation tasks  
34 to enable generalization, rather than being collected from narrow domains or a closed set of object  
35 categories. **(II)** The model needs to accurately produce instruction-conditioned segmentation masks:  
36 not only locating where robots can interact with, but also adapting to the instruction, since the same  
37 object affords different regions under different tasks. **(III)** To make the interaction actionable for  
38 robots, the model must further predict *how* the interaction should be performed, with a 3D motion

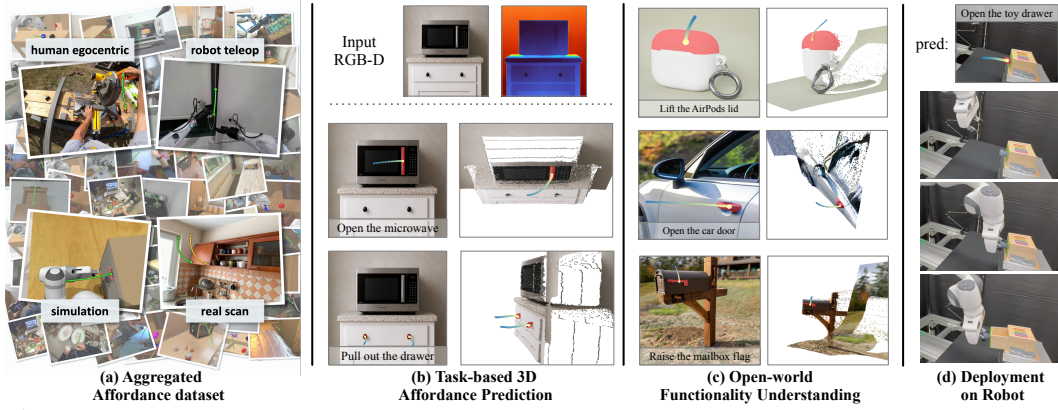


Figure 1: **Overview of AFUN.** (a) We first build a data pipeline to gather a large-scale diverse dataset for affordance understanding. (b) With such a dataset, we then train AFUN to predict a task-conditional functional segmentation mask and a 3D motion trajectory, conditioned on an RGB-D observation and a language task phrase. (c) AFUN can generalize to open-world images for functionality understanding and is directly deployable to the real robot for manipulation (d).

39 representation that a robot can follow. The 3D motion should remain expressive enough to capture  
 40 diverse behaviors and structured enough for stable supervision and robot execution.

41 In practice, however, existing affordance methods focus mainly on the second requirement alone,  
 42 formulating the problem as static segmentation [72, 46], keypoint detection [87], or reasoning-based  
 43 grounding [69]. These approaches can localize interaction regions, but they do not characterize  
 44 how the object should move after interaction. For the methods focusing on the third requirement,  
 45 some predict motion in 2D [78, 2], leaving robot execution ambiguous when lifting into 3D, while  
 46 others [86] require heuristic localization of actionable objects. Beyond limitations on each modality,  
 47 most current deep-learning models for affordance understanding [78, 2, 7, 86] still fall short of  
 48 open-world generalization due to small-scale datasets with limited diversity.

49 To address these gaps, we present **AFUN**, a step toward an open-world affordance foundation  
 50 model. First, we build a large-scale standardized data pipeline that converts public robot, human,  
 51 and simulation datasets into coherent affordance data with task descriptions, functional masks and  
 52 3D motions, extending the current affordance dataset towards one of the largest public affordance  
 53 datasets to date (Figure 1 (a)). Then, we introduce a unified affordance foundation model that jointly  
 54 predicts *where* to act (functional segmentation) and *how* the interaction should happen (3D motion,  
 55 represented as a Bézier spline curve), conditioned on text instructions from users and robot-native  
 56 RGB-D observations, as shown in Figure 1 (b). The mask can then be unprojected into 3D points for  
 57 robots to perform action leveraging downstream grasping modules such as AnyGrasp [18].

58 We evaluate AFUN on 8 segmentation-based affordance benchmarks and 3 motion-based benchmarks.  
 59 AFUN reaches **69.3** mean segmentation gIoU, compared with **45.4** for the strongest segmentation  
 60 baseline [69]; for motion prediction, it surpasses standalone baselines [78, 2, 7] by a substantial  
 61 margin. AFUN also demonstrates strong generalization capability when qualitatively evaluated on  
 62 open-world images (Fig. 1(c)). Furthermore, we deploy AFUN on a real robot for manipulation.  
 63 Without any robot-specific finetuning, AFUN can predict precise mask and motion for robot to plan  
 64 and execute a successful path for manipulation, as illustrated in Fig. 1 (d) and Fig. 7.

## 65 2 Related Work

66 **Affordance localization.** Affordance localization mainly asks *where* a task-specified interaction  
 67 is possible, differing by grounding representation. Dense 2D methods cover classical instance- and  
 68 part-level segmentation [54, 14], weakly supervised cross-view grounding [45, 37, 29, 77], egocentric-  
 69 and human-video mask supervision [38, 24, 46], LISA-style SAM grounding from LLM hidden  
 70 states [36, 57], two-stage VLM-to-segmer pipelines using LLM-emitted coordinates or boxes [69,  
 71 40, 26, 6], and language-conditioned SAM-style benchmarks and decoders [72, 67, 30, 27]. Sparse  
 72 alternatives predict contact points or keypoints, including language-conditioned image keypoints [87],  
 73 3D keypoint quadruplets encoding contact location and direction [43], and cross-category contact  
 74 transfer by semantic correspondence or retrieval [32, 35]. 3D approaches use point clouds, from  
 75 foundational benchmarks and 2D-to-3D interaction grounding [13, 79] to open-vocabulary and  
 76 language-conditioned 3D-MLLM grounding [55, 59, 44, 92, 9, 84, 70], cross-instance or object-to-  
 77 object transfer [66, 15], and video-driven MLLM learning from human-object interaction [68, 47].  
 78 Hand-centric work localizes functional grasps and dexterous contacts, from category-level grasp

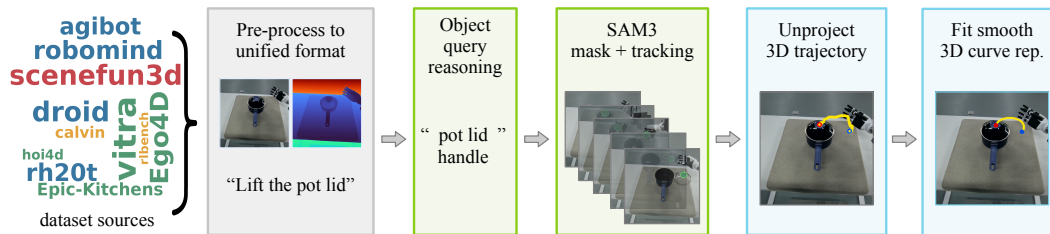


Figure 2: **Unified data collection pipeline.** We first aggregate data from various sources into a unified format (gray), then use Qwen3-VL [3] and SAM3 [49] to generate a functional affordance mask and a 2D tracking (green), and finally back-project them to obtain a 3D trajectory, which can be fit to a Bézier spline curve (blue). This scalable pipeline yields standardized, high-fidelity annotations with RGB-D observation, text phrase, mask, and 3D motion curve for training. (Best viewed in color.)

79 generation [10] and language-guided task-oriented grasping [64, 90, 85] to dexterous and finger-  
 80 specific affordance prediction [71, 22]. These lines leave post-contact motion—*how* the object should  
 81 move—largely unspecified, motivating AFUN’s joint mask-and-motion formulation.

82 **Motion representations for affordance.** Motion-focused affordance work asks *how* the object  
 83 should move after contact, with representations varying in granularity and structure. Some meth-  
 84 ods use discrete prompts or parametric articulation: elementary push/pull types tied to actionable  
 85 regions [51], scene-level functional categories with motion type and axis [12], or openable-part  
 86 motion-parameter regression for articulated objects [31, 60, 39]. Others predict continuous interaction  
 87 geometry, including dense visual action trajectories and per-point 3D articulation flow [74, 16, 89, 58],  
 88 egocentric contact heatmaps and 6-DoF object trajectories [82, 83], hand and wrist trajectories dis-  
 89 tilled from in-the-wild videos [2, 7], and 2D point tracks or 3D object-point flow as foundation  
 90 affordance [5, 86]. A third group uses affordance to condition policies, including diffusion-policy  
 91 and flow-matching action generators guided by 3D contact and post-contact trajectories [75, 91, 88],  
 92 hierarchical spatial-affordance plus low-level execution schemes [78, 53], and semantic 3D flow for  
 93 generative control [8]. Rather than outputting discrete motion types, dense flow, hand/end-effector  
 94 trajectories, or policy-conditioning signals, AFUN predicts a compact, object-centric 3D motion  
 95 curve jointly with the functional mask.

96 **Affordance data pipelines and datasets.** Affordance supervision can be categorized by annotation  
 97 target, motion source, and labeling cost. Directly annotated datasets cover early RGB-D part-  
 98 affordance benchmarks [52], image-level functional grounding datasets [45, 72, 69, 67], scene-  
 99 and shape-level 3D functional annotations [13, 12], robot-manipulation benchmarks [63, 23], and  
 100 human-video-derived corpora [46, 24, 82]. Beyond static affordance labels, trajectory and motion  
 101 supervision is obtained from internet-video hand and wrist trajectories [2, 7], egocentric 6-DoF  
 102 object trajectories with action descriptions [82, 83, 81], hand-object pose datasets used as trajectory  
 103 context [4], scene-level 3D motion benchmarks [12], and simulated articulated-object interactions [51,  
 104 74, 16, 89]. To reduce labeling cost, automatic or weakly supervised pipelines derive affordance  
 105 labels from foundation-model distillation without dense annotation [62], egocentric-video affordance  
 106 extraction [38, 24], large-scale human-behavior mining [46], part-prior weak supervision [77],  
 107 generative-AI augmentation for VLM affordance learning [23], and MLLM-assisted grounding from  
 108 human-object-interaction videos [68, 47]. Despite these efforts, existing resources remain limited in  
 109 scale, especially for 3D motion supervision; AFUN addresses this gap with an extensible pipeline  
 110 that aggregates one of the largest motion-affordance datasets to date.

### 111 3 Data Pipeline for AFUN

112 Open-world affordance learning requires a large-scale dataset that covers diverse scenarios, tasks,  
 113 objects, and action sequences while providing ground truth on *what* to manipulate (segmentation)  
 114 and *how* to manipulate (motion). Existing datasets are either too small or only contain part of the  
 115 information. In this paper, we build a unified data pipeline (Figure 2) and curate a wide range of  
 116 publicly available data, including robot demonstrations, egocentric human videos, and simulated  
 117 interactions. We annotate all the data with a common affordance schema. Each data sample contains  
 118 an RGB-D observation with a task description, a functional affordance mask, and a compact 3D  
 119 motion trajectory.

120 **Dataset Curation.** We curate datasets whose videos capture object interactions for functional  
 121 purposes and have visible action regions and object motions. Based on these criteria, we gathered  
 122 321,190 videos from 10 public sources, spanning human demonstrations [50, 21, 17, 11], robot  
 123 demonstrations [33, 17, 1, 73, 25], simulation data [48, 28], and real-world scans [12]. Since a  
 124 recording may contain multiple actions, we split each one into action intervals, resulting in 1,242,740

125 intervals to start with. This broad source data pool provides us with diverse object categories, camera  
 126 viewpoints, interaction tasks, and embodiments. Further details are provided in Appendix B.

127 **Dataset Preprocessing.** To annotate at scale, we first convert all source datasets into a common  
 128 action-interval format. Each dataset is handled with its own source-specific procedure, but the result  
 129 is the same: a set of action intervals, each paired with an observation RGB-D frame, task language,  
 130 camera pose, and the corresponding video span. We also use monocular depth estimators [61, 41] to  
 131 improve depth quality. Further details are in Appendix B.1.

132 **Annotating Object Tracks and Masks.** Prior  
 133 works often use hand or gripper trajectories as  
 134 the motion signal for affordance. However, this  
 135 can entangle the affordance-relevant post-contact  
 136 object motion with undesired pre-contact hand  
 137 motion, as shown in Figure 3 (right). Instead, we use  
 138 the tracking of the object as the post-contact  
 139 motion in our affordance foundation model (Figure 3,  
 140 left), as it directly indicates how the object moves  
 141 after contact from a robot or a human. To obtain  
 142 the tracking, we first use a vision-language model  
 143 to generate a short manipulable-part query from  
 144 the task instruction and the observation/contact  
 145 frames, then use SAM3 [49] to track the manipulated object across the action interval. This step  
 146 produces an object-centric motion trajectory and a functional affordance mask.



Figure 3: **Object trajectory vs. gripper heuristics.** Prior datasets often use hand or gripper trajectories as motion heuristics, but these can involve unwanted pre-contact motion (right). We track the object motion itself, which is more straightforward (left).

147 **Optimizing 3D Motion Curves.** With the object mask and tracking trajectory, we recover the  
 148 object’s 3D motion trajectory by back-projecting the tracked masks and taking the mean of the 3D  
 149 points in each frame. The resulting discrete path, however, is typically non-uniformly sampled and  
 150 exhibits noise due to depth estimation errors and tracking inconsistencies. To address this, we fit a  
 151 smooth parametric curve and convert it into our final canonical motion representation (Bézier spline  
 152 curve) for training. The details of curve fitting are provided in Appendix B.2.

153 **Filtering and Dataset Statistics.** Before filtering, our data pool contains 1,242,740 action intervals  
 154 spanning robot teleoperation, human egocentric recordings, simulation, and real-world scans. At  
 155 each step of processing, we filter low-quality clips, such as those with poor task grounding, occlusion,  
 156 unreliable segmentation, and insufficient motion. Each step removes around 1/2 of the samples,  
 157 eventually resulting in 223,334 samples with valid motion labels. We then perform manual annotation  
 158 and quality control, and retain 59,867 training samples for AFUN. A dataset at this scale exposes the  
 159 model to diverse interaction types, object categories, camera viewpoints, and embodiments.

## 160 4 Method

161 AFUN takes an RGB-D observation and a task phrase as input and jointly predicts a task-conditioned  
 162 functional segmentation mask, along with a 3D post-contact motion curve in a single forward pass. As  
 163 shown in Fig. 4, our model uses a simple architecture with two main components. First, we leverage  
 164 the MetaQuery mechanism [56] to connect a frozen Vision-Language Model (Qwen3-VL [3]) with  
 165 a segmentation model (SAM3 [49]) to predict functional masks. Second, we use a 3D feature  
 166 encoder [76] and a transformer decoder to predict 3D post-contact motion, represented as a Bézier  
 167 spline curve. We describe the detailed network architecture in §4.1 and the training scheme in §4.2.

### 168 4.1 Network Architecture

169 **MetaQuery Conditioning.** Introduced by Pan et al. [56], MetaQuery serves as an interface to  
 170 connect a frozen VLM with a downstream model. In particular, a small set of learnable special tokens  
 171 is appended to the VLM’s input prompt and processed through the transformer. The hidden states  
 172 in the final layer of the VLM serve as a compact conditioning feature for the downstream model.  
 173 Pan et al. [56] show that this approach can extract detailed visual conditions and transfer reasoning  
 174 capabilities to multimodal generation tasks, such as image editing.

175 We bring the MetaQuery approach to our affordance prediction model, incorporating the reason-  
 176 ing capabilities from the VLM for functional mask segmentation and motion understand-  
 177 ing. Specifically, we maintain two sets of learnable tokens:  $\mathbf{mq}^s = \{\langle \text{mq}_0^s \rangle, \dots, \langle \text{mq}_{N_s-1}^s \rangle\}$ ,  
 178  $\mathbf{mq}^m = \{\langle \text{mq}_0^m \rangle, \dots, \langle \text{mq}_{N_m-1}^m \rangle\}$ , where the first set  $\mathbf{mq}^s$  connects the VLM with the segmen-  
 179 tation model, and the second  $\mathbf{mq}^m$  connects it with the motion prediction model. The two sets of  
 180 learnable tokens are appended to the input prompt together and processed through the transformer

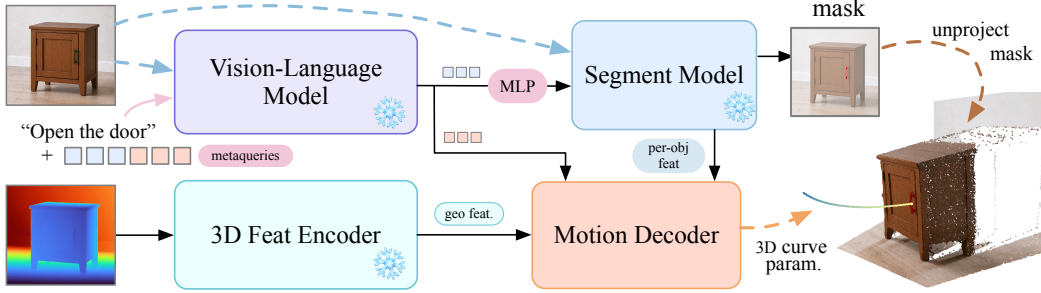


Figure 4: **AFUN architecture.** Starting from an RGB-D input and a task prompt, a frozen Qwen VLM encodes the language instruction into *semantic tokens* and *motion tokens*, and a 3D encoder converts the depth observation into geometric features. With the language information encoded, the segmentation model generates the affordance segmentation mask from the RGB, while the motion decoder takes 3D features, task-conditioned context, and per-object features to produce a relative 3D motion prediction. Together, the mask and trajectory form the final deployable 3D affordance prediction. (Best view in color.)

181 in Qwen3-VL [3]. The last hidden states of each set of tokens are then fed into the downstream  
 182 segmentation and motion model, respectively. This joint formulation allows both the segmentation  
 183 and motion models to share reasoning capabilities from VLMs within a single forward pass.

184 **Segmentation and Motion Decoding.** With the MetaQuery tokens from VLMs, we predict the  
 185 functional segmentation mask and 3D post-contact motion. For segmentation prediction, we primarily  
 186 use SAM3 [49]. Specifically, the semantic MetaQuery tokens  $\mathbf{mq}^s$  are first mapped by a two-  
 187 layer MLP into SAM3’s language-feature space. They are then passed through SAM3’s mask  
 188 decoder, which predicts per-detection boxes, masks, and object query features that are used for  
 189 motion prediction. By leveraging pretrained Qwen3-VL and SAM3, our model inherits the prior  
 190 knowledge learned from large-scale pretraining for functional segmentation understanding. For  
 191 motion prediction, we additionally encode the point cloud (unprojected from the depth input) with a  
 192 pretrained Sonata [76] network to provide 3D information. Then we use a motion decoder, which  
 193 is a transformer decoder with self-attention to the encoded 3D features and cross-attention to the  
 194 per-object features from SAM3 and the motion MetaQuery tokens  $\mathbf{mq}^m$ , to predict the parameters of  
 195 motion curves below.

196 **Curved Motion Representation.** A motion representation for open-world affordance must be  
 197 expressive enough for complex interactions yet structured enough for robust manipulation. We  
 198 therefore represent post-contact motion as an anchored 3D Bézier spline curve, parameterized by  
 199 control points. The centroid of the masked depth map defines the start point  $\mathbf{P}_0$ , and the motion  
 200 decoder predicts the remaining  $K$  ordered control points  $\{\mathbf{P}_k\}_{k=1}^K$  in relative 3D coordinates. The  
 201 trajectory is then computed with the Bernstein polynomial basis:

$$\mathbf{B}(t) = \sum_{k=0}^K \binom{K}{k} (1-t)^{K-k} t^k \mathbf{P}_k, \quad t \in [0, 1], \quad (1)$$

202 where  $\mathbf{B}(t)$  is the 3D position at normalized time  $t$ . The starting point  $\mathbf{P}_0$  anchors the curve at  
 203 the contact centroid, while the predicted control points parameterize the overall shape of the curve.  
 204 Uniformly sampling  $t \in [0, 1]$  produces executable 3D waypoints for robots.

## 205 4.2 Training Scheme

206 Directly training the full model is unstable: randomly initialized MetaQuery tokens provide a poor  
 207 conditioning signal for SAM3, and noisy mask predictions would in turn make motion supervision  
 208 ambiguous. We therefore train our model in three stages: (I) aligning the MetaQuery interface  
 209 with SAM3, (II) learning reliable task-conditioned affordance segmentation, and (III) fine-tuning  
 210 motion prediction when the model is already robust in segmentation prediction. The pretrained priors,  
 211 Qwen-VL, SAM3, and Sonata, are kept frozen throughout the training.

212 **Stage 1: MetaQuery–SAM3 Alignment.** Prior to end-to-end training, we initialize and train the  
 213 MetaQuery tokens and projection MLP by aligning Qwen-derived features with SAM3’s native text-  
 214 conditioning space on the Visual Genome dataset [34]. For each caption-image pair, we encode the  
 215 caption with SAM3’s text encoder; in parallel, Qwen3-VL processes the same caption and image, and  
 216 the projection MLP projects the resulting MetaQuery features into SAM3 text space. We then run the  
 217 SAM3 decoder with cross-attention to both the projected MetaQuery features and the original SAM3  
 218 text features. The decoder hidden states from the two branches are optimized with a Mean-Squared  
 219 Error (MSE) loss, which provides a more stable initialization than training the new tokens directly

220 from mask supervision. This alignment step yields a strong initialization for the MetaQuery tokens  
 221 and the Qwen-to-SAM3 MLP, thereby stabilizing subsequent joint affordance training.

222 **Stage 2: End-to-End Training for Affordance Segmentation.** In the second stage, we train our  
 223 affordance segmentation model end-to-end on an aggregated mixture of four affordance datasets:  
 224 HOVA-500K [46], RAGNet [72], InstructPart [67], and ReasonAFF [69]. The unfrozen parameters  
 225 are identical to those in Stage 1: the MetaQuery tokens and the projection MLP. The motion prediction  
 226 branch is disabled, and we only train the model with objectives from SAM3 [49], which combines  
 227 Hungarian-matched box regression ( $\ell_1 + \text{GIoU}$ ), presence classification, per-query mask prediction  
 228 (focal BCE + Dice), and a semantic-segmentation term (focal + Dice + presence), all averaged with  
 229 the same hyperparameters as SAM3. We refer readers to the original paper for more details.

230 **Stage 3: Joint Motion and Segmentation Training.** In the final stage, we train segmentation  
 231 and motion prediction jointly on our own aggregated affordance dataset curated from Section 3,  
 232 together with the Stage 2 training data. The total objective combines the Stage 2 SAM3 grounding  
 233 loss  $\mathcal{L}_{\text{sam3}}$ , down-weighted to prevent the segmentation head from overfitting, with a curve loss  
 234  $\mathcal{L}_{\text{curve}}$  on sampled trajectory points to learn the motion:

$$\mathcal{L} = \lambda_{\text{sam3}} \mathcal{L}_{\text{sam3}} + \lambda_{\text{curve}} \mathcal{L}_{\text{curve}}. \quad (2)$$

235 We follow the point-sampling loss from Curve-GCN [42] to supervise motion prediction. Specifically,  
 236 for each SAM3-matched query  $(b, q) \in \mathcal{M}$  returned by the Hungarian matcher, we evaluate both the  
 237 predicted Bézier curve  $\widehat{\mathbf{B}}_{b,q}(t)$  and its matched ground-truth curve  $\mathbf{B}_{b,q}^*(t)$  on a fixed uniform time  
 238 interval  $\{t_i = i/(T - 1)\}_{i=0}^{T-1}$  and minimize the  $\ell_1$  distance between the sampled points,

$$\mathcal{L}_{\text{curve}} = \frac{1}{|\mathcal{M}|T} \sum_{(b,q) \in \mathcal{M}} \sum_{i=0}^{T-1} \|\widehat{\mathbf{B}}_{b,q}(t_i) - \mathbf{B}_{b,q}^*(t_i)\|_1. \quad (3)$$

239 In practice, we find this point-sampling supervision substantially more effective than directly regress-  
 240 ing the locations of control points.

## 241 5 Experiments

### 242 5.1 Implementation Details

243 We use Qwen3-VL-8B [3] as our VLM backbone, SAM3 [49] as the segmentation model, and  
 244 Sonata [76] as the 3D feature encoder; all three pretrained components are frozen throughout training.  
 245 The motion decoder uses six transformer layers. Along with the MLPs and MetaQuery tokens, our  
 246 model adds only **32.21M** trainable parameters on top of the pretrained models. We use 64 MetaQuery  
 247 tokens in total, where each semantic and motion branch has 32. We set  $\lambda_{\text{SAM}} = 0.5$ ,  $\lambda_{\text{curve}} = 100$ ,  
 248 and train the model with a learning rate of  $2 \times 10^{-4}$ . For the point-sampling loss, we sample  $T = 16$   
 249 points per curve. We train AFUN on  $4 \times \text{NVIDIA GH200}$  GPUs for approximately eight days. The  
 250 three stages mentioned above use batch sizes of 196, 128, and 96, respectively, and run for 10,000,  
 251 40,000, and 20,000 steps, respectively.

### 252 5.2 Affordance Evaluation

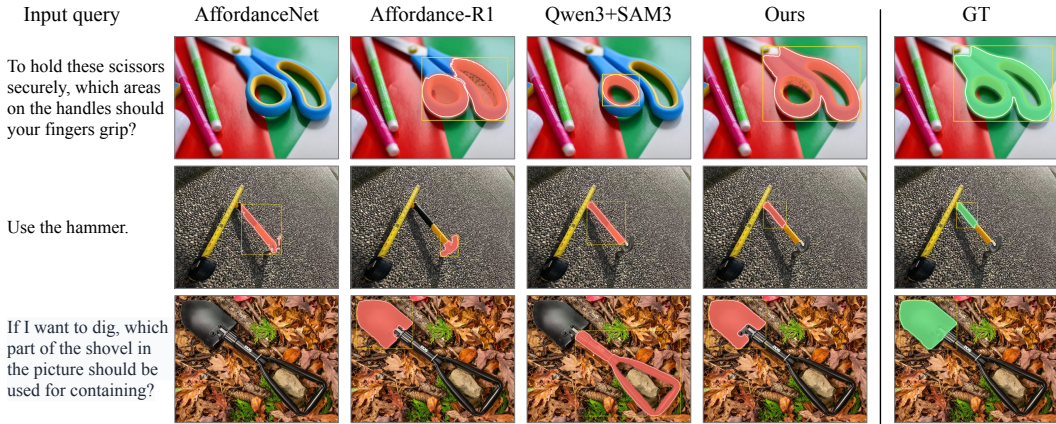
253 To comprehensively demonstrate the affordance understanding capability of AFUN, we evaluate  
 254 it from three perspectives: accuracy of affordance mask segmentation, quality of the contact point  
 255 derived from the mask, and quality of 3D motion.

#### 256 5.2.1 Affordance Segmentation Evaluation

257 We first evaluate AFUN’s capability to reason about *where* the affordance lies by measuring seg-  
 258 mentation quality. We compare against three baselines: a zero-shot Qwen3-VL-8B [3] object query  
 259 generation + SAM3 [49] mask generation pipeline, AffordanceNet [72], and Affordance-R1 [69].

260 We show qualitative task-conditioned affordance mask segmentation results in Fig. 5. AFUN  
 261 consistently predicts the correct affordance region for the diverse task instructions, can precisely  
 262 segment complex regions (scissors handle with holes), and strictly aligns with the given task (shovel  
 263 blade-containing, hammer handle-using), demonstrating superior performance in reasoning about  
 264 the task-specific affordance compared to baselines.

265 Quantitatively, following our baselines [72, 46, 69], we evaluate on eight test sets drawn from four  
 266 affordance benchmarks, and report gIoU and cIoU metrics in Table 1. Across all test sets, AFUN  
 267 outperforms all the baselines, achieves the best gIoU and cIoU, and improves overall mean gIoU/cIoU  
 268 by 23.9/26.3 points over the strongest baseline. Notably, even when using the Qwen3-VL-2B variant  
 269 with fewer parameters, AFUN remains superior to the baseline models by a large margin.



**Figure 5: Qualitative Examples on Affordance Segmentation.** AFUN accurately segments task-specific affordance regions, including complex scissor handles with holes and intent-dependent regions such as shovel blades for containing or hammer handles for using. We provide more examples in Appendix B.3. Quantitative results on affordance mask segmentation (IoU / cIoU, %; higher is better). Best results are **bolded**: AFUN significantly outperforms all the baselines across 8 datasets in both metrics. Even with a smaller 2B model, AFUN still improves over the baselines by a large margin.

Method	HANDAL-Mini	3DOI	HANDAL Easy	HANDAL Hard	3DOI Easy	HOVA-500K	ReasonAFF	InstructPart	Mean
Qwen3-VL-8B [3] + SAM3 [49]	40.3 / 44.4	36.4 / 22.5	44.8 / 50.6	45.2 / 52.2	44.3 / 27.2	63.3 / 38.5	31.9 / 19.6	41.4 / 28.1	42.5 / 36.5
Affordance-R1 [69]	18.2 / 10.6	51.8 / 38.5	39.8 / 32.0	37.4 / 26.5	61.4 / 53.6	27.9 / 14.1	67.1 / 62.1	67.2 / 58.6	45.4 / 36.2
AffordanceNet [72]	59.6 / 58.8	39.9 / 37.8	55.2 / 50.7	52.8 / 50.3	36.3 / 36.6	51.7 / 28.1	25.4 / 19.9	29.0 / 21.8	45.0 / 40.9

Query	Ours	VRB	VidBot	A0	General-Flow <sup>†</sup>	Ours	GT
Put the lid on the pot.							
Shut the door of the orange microwave on the desktop.							
Close the door.							

**Figure 6: Qualitative motion prediction results.** AFUN accurately localizes the actionable object region and predicts smooth, task-aligned 3D motion curves, whereas the baselines often fail to identify the relevant affordance region or produce physically plausible motion. <sup>†</sup> General Flow used the mask prediction from our AFUN for its starting query points.

## 270 5.2.2 Contact Point Evaluation

271 Beyond using masks for affordance, prior work also adopts contact points as an affordance representation; we therefore compare with A0 [78], GLOVER++ [46], VRB [2], and measure whether a  
 272 predicted contact point lies on the ground-truth affordance mask. For AFUN, we take the Pole of  
 273 Inaccessibility [19] of the predicted mask as the contact point. We use hit rate  $\Pr[\text{point} \in \text{GT mask}]$ ,  
 274 which measures whether the predicted contact point lies on the affordance mask as the evaluation  
 275 metric. As shown in Table 2, AFUN significantly outperforms the best baseline by 12.7%–61.3%  
 276 (55.7% on InstructPart and 61.3% on ReasonAFF).  
 277

## 278 5.2.3 3D Motion Evaluation

279 **Evaluation Datasets.** We evaluate 3D motion on three test sets with different domain shifts. (I)  
 280 The AFUN test set (121 examples) is a cross-source split randomly sampled from the high-quality  
 281 set we curated. This test set is further verified through a second human quality-control pass, and

282 we exclude these data samples from the training set to prevent data leak. (II) The SceneFun3D [12]  
 283 test set (721 examples) comes from the original validation set in SceneFun3D and contains scenes  
 284 that are not present in the training set. For each task in each scene, we use the first frame in which  
 285 the target object is visible for evaluation (details of dataset processing in Appendix B.1). (III) The  
 286 RoboMIND2 dataset test set (156 examples) is an out-of-domain test set deliberately excluded from  
 287 training. We keep functionality-related tasks and remove relocation-only instructions such as “place  
 288 A to B” for evaluation, as such waypoints are usually non-deterministic.

289 **Evaluation Metrics and Baselines.** We evaluate predicted 3D motion curves using Average Dis-  
 290 placement Error (ADE), Final Displacement Error (FDE) computed in both absolute scale and relative  
 291 scale, contact-in-mask hit rate (CIM). We compare AFUN with four 3D affordance baselines: A0 [78],  
 292 VRB [2], VidBot [7], and General Flow [86]. For each baseline, we follow their official protocol to  
 293 obtain the 3D motion predictions, and linearly interpolate every prediction and every ground-truth  
 294 trajectory to a common length of  $T=50$  points for evaluation. Note that General Flow [86] requires a  
 user-picked starting point, and we use the predicted mask from our model as the starting point for it.

Table 2: Contact point evaluation with point hit rate (%; higher is better). We compare with 2D point-based methods. Best per column in **bold**. We use the Pole of Inaccessibility of the predicted mask as the predicted point.

Method	HANDAL-Mini	3DOI	HANDAL Easy	HANDAL Hard	3DOI Easy	HOVA-500K	ReasonAFF	InstructPart
A0 [78]	4.6	3.6	5.0	6.0	3.6	3.8	20.7	22.7
VRB [2]	31.5	46.6	31.5	31.3	48.4	22.4	35.2	39.8
GLOVER++ [46]	67.6	6.8	39.2	34.4	4.9	28.6	4.3	4.7
AFUN (Ours)	<b>80.3</b>	<b>67.3</b>	<b>88.2</b>	<b>88.8</b>	<b>82.6</b>	<b>88.3</b>	<b>96.5</b>	<b>95.5</b>

296 **Evaluation Results.** We provide quantitative results in Table 3 and qualitative results in Fig. 6.  
 297 AFUN achieves the best ADE and FDE in both absolute and relative scale on all three test sets,  
 298 and significantly outperforms the baselines in CIM. Even when General Flow is evaluated under a  
 299 favorable protocol that provides it with AFUN’s predicted mask and start anchor, AFUN still achieves  
 300 substantially better motion prediction results. This advantage is further shown in Fig. 6: AFUN  
 301 produces task-aligned masks and motions, whereas the baselines often produce both implausible  
 302 object localization and task-inconsistent trajectories.

### 303 5.3 Ablations

304 We ablate three different design choices of our model: the LLM backbone, the 3D feature encoder,  
 305 and the motion curve parameterization. Results are provided in Table 4 and Table 5.

306 For different LLM backbones, we train the model using the same recipe as our default model and  
 307 report the evaluation performance on all the 8 test sets. Our default model with Qwen3-VL-8B  
 308 achieves the best segmentation performance, outperforming both the smaller model Qwen3-VL-2B  
 309 and the larger Qwen3.5-9B. We hypothesize: the reason why larger Qwen3.5-9B underperforms is  
 310 that its general-purpose MoE design might be less suited to dense vision–language prediction.

311 For the 3D feature encoder, we replace Sonata with the *emphDFormerv2* [80] architecture and train it  
 312 using the same recipe. We evaluate the performance on the open-domain RoboMind2 test set. Our  
 313 default 3D feature encoder outperforms *DFormerv2* [80], benefiting from stronger 3D geometric cues  
 314 in point cloud-derived features.

315 For motion curve representation, we compare our curve parameterization with the representation  
 316 in *OPD* [31]. Our representation achieves better results, as the single parameterization for multiple  
 317 motion types in *OPD* can introduce ambiguity.

### 318 5.4 Real-Robot Demonstration

319 Deploying AFUN on real robotic platforms is straightforward and requires no additional task-  
 320 specific heuristics. Given a calibrated RGB-D input from one camera, AFUN predicts a contact  
 321 mask and post-contact motion trajectory; the mask is back-projected to localize the target object,  
 322 while AnyGrasp [18] estimates feasible grasp poses from the reconstructed scene point cloud. The  
 323 predicted trajectory, represented as a smooth spline curve, provides a local tangent direction for  
 324 adapting the gripper orientation, enabling rotational manipulation such as opening a microwave.  
 325 This orientation-aware execution is difficult to obtain from line-based trajectory predictions in prior  
 326 approaches [2, 78, 86].

**Table 3:** Quantitative 3D motion evaluation. ADE/FDE are in meters; subscript  $a$  is absolute,  $r$  is relative. CIM is the contact-in-mask hit rate. Best per dataset in **bold**. General Flow<sup>†</sup> gives no starting point  $\mathbf{r}_0$  for motion prediction, and uses predicted mask from our AFUN to get its query points. Yet, it still underperforms our model.

Dataset	Method	ADE <sub>a</sub> ↓	FDE <sub>a</sub> ↓	ADE <sub>r</sub> ↓	FDE <sub>r</sub> ↓	CIM %↑	#fail↓
AFUN test set ( $n=121$ )	A0 [78]	0.378	0.369	0.140	0.284	6.6	<b>0</b>
	VRB [2]	0.242	0.297	0.087	0.220	11.0	3
	VidBot [7]	0.520	0.614	0.192	0.364	6.7	2
	General Flow <sup>†</sup> [86]	0.110	0.230	0.088	0.225	–	<b>0</b>
	AFUN (Ours)	<b>0.098</b>	<b>0.139</b>	<b>0.080</b>	<b>0.135</b>	<b>81.0</b>	<b>0</b>
SceneFun3D test ( $n=721$ )	A0 [78]	1.008	1.066	0.249	0.476	3.3	<b>0</b>
	VRB [2]	0.606	0.702	0.227	0.446	11.8	36
	VidBot [7]	0.772	0.848	0.201	0.393	9.4	11
	General Flow <sup>†</sup> [86]	0.413	0.572	0.212	0.413	–	4
	AFUN (Ours)	<b>0.351</b>	<b>0.441</b>	<b>0.135</b>	<b>0.260</b>	<b>67.3</b>	1
RoboMIND2 test ( $n=156$ )	A0 [78]	0.374	0.388	0.193	0.327	3.2	<b>0</b>
	VRB [2]	0.299	0.373	0.190	0.330	27.4	10
	VidBot [7]	0.368	0.479	0.240	0.414	23.7	5
	General Flow <sup>†</sup> [86]	0.260	0.369	0.184	0.314	–	2
	AFUN (Ours)	<b>0.254</b>	<b>0.323</b>	<b>0.177</b>	<b>0.276</b>	<b>62.2</b>	<b>0</b>

**Table 4:** LLM backbone ablation.

Variant	Mean gIoU↑	Mean cIoU↑
AFUN (Ours)	<b>69.6</b>	<b>66.4</b>
w/ Qwen3-VL-2B	66.6	61.8
w/ Qwen3.5-9B	61.8	54.4

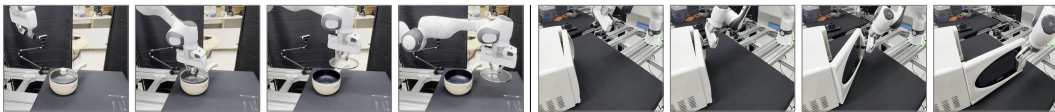
**Table 5:** 3D motion ablations on RoboMIND2.

Variant	ADE <sub>a</sub> ↓	FDE <sub>a</sub> ↓
AFUN (Ours)	<b>0.254</b>	<b>0.323</b>
w/ DFormerv2 [80] (3D feat. encoder)	0.273	0.348
w/ OPD [31] (curve parameterization)	0.282	0.374

327 We evaluate AFUN on four real-world tasks: Pick Up Screwdriver,  
 328 Take Off Pot Lid, Open Drawer, and Open Microwave, using a  
 329 Franka Research 3 arm and two calibrated third-person RGB-D  
 330 RealSense D435 cameras. For each task, AFUN uses one RGB-D  
 331 observation as input, while observations from both cameras are  
 332 fused into the scene point cloud used by AnyGrasp. We report  
 333 success rates in Tab. 6 and qualitative examples in Fig. 7. AFUN achieves an average success  
 334 rate of 90%, demonstrating reliable real-robot deployment for both contact-centric grasping and  
 335 orientation-aware articulated-object manipulation.

**Table 6:** Real-world Task Performance.

Task	Success Rate
Pick Up Screwdriver	1.0
Take Off Pot Lid	1.0
Open Drawer	0.8
Open Microwave	0.8



**Figure 7: Real-robot deployment (Franka).** AFUN can be directly deployed to a real robotic system without any additional task-specific heuristics. Given a task from the user, our model can accurately locate the actionable (grasping) region and produce an accurate post-contact trajectory for robot manipulation.

## 336 6 Conclusion

337 In this paper, we present **AFUN**, a step towards an affordance foundation model for functionality  
 338 understanding. From a single RGB-D observation and a language task description, AFUN predicts  
 339 a task-conditional functional mask (*where* to interact) and a 3D post-contact motion curve (*how*  
 340 to interact). To achieve open-world generalization, we build a large-scale standardized data pipeline that  
 341 converts heterogeneous robot, human, simulation, and real-world scan data into a shared affordance  
 342 schema with language, masks, and object-centric 3D motion annotations. Empirically, AFUN  
 343 outperforms all baselines on affordance segmentation across eight test sets from four benchmarks;  
 344 predicts substantially more accurate contact points; and achieves the best 3D motion performance on  
 345 all three motion test sets. Without embodiment-specific finetuning, AFUN can be directly deployed  
 346 in the real robot for manipulation, suggesting a practical path towards open-world affordance models  
 347 that unify functionality perception with executable action. We provide limitations, failure cases, and  
 348 future directions in Appendix 6.

## References

- 349
- 350 [1] AgiBot-World-Contributors, Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan  
351 Feng, Shenyuan Gao, Xindong He, Xu Huang, et al. AgiBot World Colosseo: A large-  
352 scale manipulation platform for scalable and intelligent embodied systems. *arXiv preprint*  
353 *arXiv:2503.06669*, 2025.
- 354 [2] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances  
355 from human videos as a versatile representation for robotics. In *Proceedings of the IEEE/CVF*  
356 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- 357 [3] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao  
358 Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie  
359 Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei Li, Kaixin  
360 Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang Liu, Dayiheng  
361 Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen Meng, Xuancheng  
362 Ren, Xingzhang Ren, Sibao Song, Yuchong Sun, Jun Tang, Jianhong Tu, Jianqiang Wan, Peng  
363 Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao Xie, Yiheng Xu, Haiyang Xu, Jin  
364 Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang, Bowen Yu, Fei Zhang, Hang Zhang,  
365 Xi Zhang, Bo Zheng, Humen Zhong, Jingren Zhou, Fan Zhou, Jing Zhou, Yuanzhi Zhu, and  
366 Ke Zhu. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025.
- 367 [4] Prithviraj Banerjee, Sindi Shkodrani, Pierre Moulon, Shreyas Hampali, Shangchen Han, Fan  
368 Zhang, Linguang Zhang, Jade Fountain, Edward Miller, Selen Basol, Richard Newcombe,  
369 Robert Wang, Jakob Julian Engel, and Tomas Hodan. HOT3D: Hand and object tracking in 3D  
370 from egocentric multi-view videos. In *Proceedings of the IEEE/CVF Conference on Computer*  
371 *Vision and Pattern Recognition (CVPR)*, pages 7061–7071, 2025.
- 372 [5] Homanga Bharadhwaj, Roozbeh Mottaghi, Abhinav Gupta, and Shubham Tulsiani. Track2Act:  
373 Predicting point tracks from internet videos enables generalizable robot manipulation. In  
374 *European Conference on Computer Vision (ECCV)*. Springer, 2024.
- 375 [6] Changmao Chen, Yuren Cong, and Zhen Kan. Worldafford: Affordance grounding based on  
376 natural language instructions. In *36th IEEE International Conference on Tools with Artificial*  
377 *Intelligence, ICTAI 2024, Herndon, VA, USA, October 28-30, 2024*, pages 822–828. IEEE, 2024.  
378 doi: 10.1109/ICTAI62512.2024.00120. URL [https://doi.org/10.1109/ICTAI62512.](https://doi.org/10.1109/ICTAI62512.2024.00120)  
379 2024.00120.
- 380 [7] Hanzhi Chen, Boyang Sun, Anran Zhang, Marc Pollefeys, and Stefan Leutenegger. VidBot:  
381 Learning generalizable 3D actions from in-the-wild 2D human videos for zero-shot robotic  
382 manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
383 *Recognition (CVPR)*, pages 27661–27672, 2025.
- 384 [8] Tianxing Chen, Yao Mu, Zhixuan Liang, Zanxin Chen, Shijia Peng, Qiangyu Chen, Mingkun  
385 Xu, Ruizhen Hu, Hongyuan Zhang, Xuelong Li, and Ping Luo. G3Flow: Generative 3D  
386 semantic flow for pose-aware and generalizable object manipulation. In *Proceedings of the*  
387 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1735–1744,  
388 June 2025.
- 389 [9] Hengshuo Chu, Xiang Deng, Qi Lv, Xiaoyang Chen, Yinchuan Li, Jianye Hao, and Liqiang Nie.  
390 3d-affordancellm: Harnessing large language models for open-vocabulary affordance detection  
391 in 3d worlds. In *The Thirteenth International Conference on Learning Representations, ICLR*  
392 *2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025. URL [https://openreview.](https://openreview.net/forum?id=GThTiuXgDC)  
393 [net/forum?id=GThTiuXgDC](https://openreview.net/forum?id=GThTiuXgDC).
- 394 [10] Enric Corona, Albert Pumarola, Guillem Alenya, Francesc Moreno-Noguer, and Grégory Rogez.  
395 GanHand: Predicting human grasp affordances in multi-object scenes. In *Proceedings of the*  
396 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5031–5041,  
397 2020.
- 398 [11] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Jian Ma, Evangelos  
399 Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray.  
400 Rescaling egocentric vision: Collection, pipeline and challenges for EPIC-KITCHENS-100.  
401 *International Journal of Computer Vision (IJCV)*, 130(1):33–55, 2022.

- 402 [12] Alexandros Delitzas, Ayça Takmaz, Federico Tombari, Robert Sumner, Marc Pollefeys, and  
403 Francis Engelmann. SceneFun3D: Fine-grained functionality and affordance understanding  
404 in 3D scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
405 Recognition (CVPR)*, 2024.
- 406 [13] Shengheng Deng, Xun Xu, Chaozheng Wu, Ke Chen, and Kui Jia. 3d affordancenet: A  
407 benchmark for visual object affordance understanding. In *Proceedings of the IEEE Conference  
408 on Computer Vision and Pattern Recognition*, 2021.
- 409 [14] Thanh-Toan Do, Anh Nguyen, and Ian Reid. Affordancenet: An end-to-end deep learning ap-  
410 proach for object affordance detection. In *International Conference on Robotics and Automation  
411 (ICRA)*, 2018.
- 412 [15] Xiaoxiang Dong and Weiming Zhi. Affordance transfer across object instances via semantically  
413 anchored functional map. *CoRR*, abs/2602.14874, 2026. doi: 10.48550/ARXIV.2602.14874.  
414 URL <https://doi.org/10.48550/arXiv.2602.14874>.
- 415 [16] Ben Eisner, Harry Zhang, and David Held. FlowBot3D: Learning 3D articulation flow to  
416 manipulate articulated objects. In *Proceedings of Robotics: Science and Systems*, New York  
417 City, NY, USA, June 2022. doi: 10.15607/RSS.2022.XVIII.018.
- 418 [17] Hao-Shu Fang, Hongjie Fang, Zhenyu Tang, Jirong Liu, Junbo Wang, Haoyi Zhu, and Cewu Lu.  
419 RH20T: A comprehensive robotic dataset for learning diverse skills in one-shot. *arXiv preprint  
420 arXiv:2307.00595*, 2023.
- 421 [18] Hao-Shu Fang, Chenxi Wang, Hongjie Fang, Minghao Gou, Jirong Liu, Hengxu Yan, Wenhai  
422 Liu, Yichen Xie, and Cewu Lu. Anygrasp: Robust and efficient grasp perception in spatial and  
423 temporal domains. *IEEE Transactions on Robotics*, 2023. doi: 10.1109/TRO.2023.3281153.
- 424 [19] Daniel Garcia-Castellanos and Umberto Lombardo. Poles of inaccessibility: A calculation  
425 algorithm for the remotest places on earth. *Scottish Geographical Journal*, 123(3):227–233,  
426 September 2007. ISSN 1751-665X. doi: 10.1080/14702540801897809. URL [http://dx.  
427 doi.org/10.1080/14702540801897809](http://dx.doi.org/10.1080/14702540801897809).
- 428 [20] James J. Gibson. *The Ecological Approach to Visual Perception*. Houghton Mifflin, 1979.
- 429 [21] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit  
430 Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4D: Around the world  
431 in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer  
432 Vision and Pattern Recognition (CVPR)*, pages 18995–19012, 2022.
- 433 [22] Yifan Han, Yichuan Peng, Pengfei Yi, Junyan Li, Hanqing Wang, Gaojing Zhang, Qi Peng  
434 Liu, and Wenzhao Lian. Fsag: Enhancing human-to-dexterous-hand finger-specific affordance  
435 grounding via diffusion models, 2026. URL <https://arxiv.org/abs/2601.08246>.
- 436 [23] Xiaoshuai Hao, Yingbo Tang, Lingfeng Zhang, Yanbiao Ma, Yunfeng Diao, Ziyu Jia, Wenbo  
437 Ding, Hangjun Ye, and Long Chen. Roboafford++: A generative ai-enhanced dataset for  
438 multimodal affordance learning in robotic manipulation and navigation. *arXiv preprint  
439 arXiv:2511.12436*, 2025.
- 440 [24] Marvin Heidinger, Snehal Jauhri, Vignesh Prasad, and Georgia Chalvatzaki. 2handedafforder:  
441 Learning precise actionable bimanual affordances from human videos. In *Proceedings of the  
442 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14743–14753, October  
443 2025.
- 444 [25] Chengkai Hou, Kun Wu, Jiaming Liu, Zhengping Che, Di Wu, Fei Liao, Guangrun Li, Jingyang  
445 He, Qiuxuan Feng, Zhao Jin, et al. RoboMIND 2.0: A multimodal, bimanual mobile ma-  
446 nipulation dataset for generalizable embodied intelligence. *arXiv preprint arXiv:2512.24653*,  
447 2025.
- 448 [26] Siyuan Huang, Iaroslav Ponomarenko, Zhengkai Jiang, Xiaoqi Li, Xiaobin Hu, Peng Gao,  
449 Hongsheng Li, and Hao Dong. Manipvqa: Injecting robotic affordance and physically grounded  
450 information into multi-modal large language models. In *IEEE/RSJ International Conference*

- 451 *on Intelligent Robots and Systems, IROS 2024, Abu Dhabi, United Arab Emirates, October*  
452 *14-18, 2024*, pages 7580–7587. IEEE, 2024. doi: 10.1109/IROS58592.2024.10801993. URL  
453 <https://doi.org/10.1109/IROS58592.2024.10801993>.
- 454 [27] Yizhou Huang, Fan Yang, Guoliang Zhu, Gen Li, Hao Shi, Yukun Zuo, Wenrui Chen, Zhiyong  
455 Li, and Kailun Yang. Resource-efficient affordance grounding with complementary depth and  
456 semantic prompts. In *IEEE/RSJ International Conference on Intelligent Robots and Systems,*  
457 *IROS 2025, Hangzhou, China, October 19-25, 2025*, pages 7788–7795. IEEE, 2025. doi:  
458 10.1109/IROS60139.2025.11245943. URL [https://doi.org/10.1109/IROS60139.2025.](https://doi.org/10.1109/IROS60139.2025.11245943)  
459 11245943.
- 460 [28] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J. Davison. RL Bench: The  
461 robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*  
462 *(RA-L)*, 5(2):3019–3026, 2020.
- 463 [29] Ji Ha Jang, Hoigi Seo, and Se Young Chun. INTRA: Interaction relationship-aware weakly  
464 supervised affordance grounding. In *European Conference on Computer Vision (ECCV)*, pages  
465 18–34. Springer, 2024.
- 466 [30] Dengyang Jiang, Zanyi Wang, Hengzhuang Li, Sizhe Dang, Teli Ma, Wei Wei, Guang Dai, Lei  
467 Zhang, and Mengmeng Wang. Affordancesam: Segment anything once more in affordance  
468 grounding, 2025. URL <https://arxiv.org/abs/2504.15650>.
- 469 [31] Hanxiao Jiang, Yongsan Mao, Manolis Savva, and Angel X. Chang. OPD: Single-view 3d  
470 openable part detection. In *European Conference on Computer Vision (ECCV)*, 2022.
- 471 [32] Yuanchen Ju, Kaizhe Hu, Guowei Zhang, Gu Zhang, Mingrun Jiang, and Huazhe Xu. Robo-  
472 abc: Affordance generalization beyond categories via semantic correspondence for robot  
473 manipulation. In *European Conference on Computer Vision*, pages 222–239. Springer, 2024.
- 474 [33] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth  
475 Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis,  
476 et al. DROID: A large-scale in-the-wild robot manipulation dataset. In *Proceedings of Robotics:*  
477 *Science and Systems (RSS)*, 2024. arXiv:2403.12945.
- 478 [34] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie  
479 Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei.  
480 Visual genome: Connecting language and vision using crowdsourced dense image annotations.  
481 *International Journal of Computer Vision (IJCV)*, 123(1):32–73, 2017.
- 482 [35] Yuxuan Kuang, Junjie Ye, Haoran Geng, Jiageng Mao, Congyue Deng, Leonidas Guibas,  
483 He Wang, and Yue Wang. RAM: Retrieval-based affordance transfer for generalizable zero-  
484 shot robotic manipulation. In *Proceedings of The 8th Conference on Robot Learning (CoRL)*,  
485 volume 270 of *Proceedings of Machine Learning Research*, pages 547–565. PMLR, 2024. URL  
486 <https://proceedings.mlr.press/v270/kuang25a.html>.
- 487 [36] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. LISA:  
488 Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference*  
489 *on Computer Vision and Pattern Recognition (CVPR)*, pages 9579–9589, 2024.
- 490 [37] Gen Li, Varun Jampani, Deqing Sun, and Laura Sevilla-Lara. Locate: Localize and transfer  
491 object parts for weakly supervised affordance grounding. In *Proceedings of the IEEE/CVF*  
492 *Conference on Computer Vision and Pattern Recognition*, 2023.
- 493 [38] Gen Li, Nikolaos Tsagkas, Jifei Song, Ruaridh Mon-Williams, Sethu Vijayakumar, Kun Shao,  
494 and Laura Sevilla-Lara. Learning precise affordances from egocentric videos for robotic  
495 manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
496 2025.
- 497 [39] Siqi Li, Xiaoxue Chen, Haoyu Cheng, Guyue Zhou, Hao Zhao, and Guanzhong Tian. Locate n’  
498 Rotate: Two-stage openable part detection with geometric foundation model priors. In Minsu  
499 Cho, Ivan Laptev, Du Tran, Angela Yao, and Hongbin Zha, editors, *Computer Vision - ACCV*  
500 *2024 - 17th Asian Conference on Computer Vision, Hanoi, Vietnam, December 8-12, 2024*,

- 501 *Proceedings, Part VII*, Lecture Notes in Computer Science, pages 716–732. Springer, 2024. doi:  
502 10.1007/978-981-96-0963-5\_6. URL [https://doi.org/10.1007/978-981-96-0963-5\\_](https://doi.org/10.1007/978-981-96-0963-5_6)  
503 6.
- 504 [40] Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang,  
505 Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for  
506 object-centric robotic manipulation. In *IEEE/CVF Conference on Computer Vision and Pat-*  
507 *tern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 18061–18070.  
508 IEEE, 2024. doi: 10.1109/CVPR52733.2024.01710. URL [https://doi.org/10.1109/](https://doi.org/10.1109/CVPR52733.2024.01710)  
509 [CVPR52733.2024.01710](https://doi.org/10.1109/CVPR52733.2024.01710).
- 510 [41] Haotong Lin, Sili Chen, Junhao Liew, Donny Y. Chen, Zhenyu Li, Guang Shi, Jiashi Feng, and  
511 Bingyi Kang. Depth anything 3: Recovering the visual space from any views. *arXiv preprint*  
512 *arXiv:2511.10647*, 2025.
- 513 [42] Huan Ling, Jun Gao, Amlan Kar, Wenzheng Chen, and Sanja Fidler. Fast interactive object  
514 annotation with Curve-GCN. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
515 *and Pattern Recognition (CVPR)*, 2019.
- 516 [43] Zhiyang Liu, Ruiteng Zhao, Lei Zhou, Chengran Yuan, Yuwei Wu, Sheng Guo, Zhengshen  
517 Zhang, Chenchen Liu, Marcelo H. Ang, and Francis E. H. Tay. 3d affordance keypoint  
518 detection for robotic manipulation. In *IEEE/RSJ International Conference on Intelligent*  
519 *Robots and Systems, IROS 2024, Abu Dhabi, United Arab Emirates, October 14-18, 2024*,  
520 pages 7528–7534. IEEE, 2024. doi: 10.1109/IROS58592.2024.10801792. URL [https://doi.org/10.1109/](https://doi.org/10.1109/IROS58592.2024.10801792)  
521 [IROS58592.2024.10801792](https://doi.org/10.1109/IROS58592.2024.10801792).
- 522 [44] Dongyue Lu, Lingdong Kong, Tianxin Huang, and Gim Hee Lee. Geal: Generalizable 3d  
523 affordance learning with cross-modal consistency. In *Proceedings of the IEEE/CVF Conference*  
524 *on Computer Vision and Pattern Recognition (CVPR)*, pages 1680–1690, June 2025.
- 525 [45] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. Learning affordance  
526 grounding from exocentric images. In *CVPR*, 2022.
- 527 [46] Teli Ma, Jia Zheng, Zifan Wang, Ziyao Gao, Jiaming Zhou, and Junwei Liang. GLOVER++:  
528 Unleashing the potential of affordance learning from human behaviors for robotic manipulation.  
529 In *Proceedings of The 9th Conference on Robot Learning (CoRL)*, volume 305 of *Proceedings of*  
530 *Machine Learning Research*, pages 3972–3994. PMLR, 2025. URL [https://proceedings.](https://proceedings.mlr.press/v305/ma25b.html)  
531 [mlr.press/v305/ma25b.html](https://proceedings.mlr.press/v305/ma25b.html).
- 532 [47] Aihua Mao, Kaihang Huang, Yong-Jin Liu, Chee Seng Chan, and Ying He. Vagnet: Grounding  
533 3d affordance from human-object interactions in videos, 2026. URL [https://arxiv.org/](https://arxiv.org/abs/2602.20608)  
534 [abs/2602.20608](https://arxiv.org/abs/2602.20608).
- 535 [48] Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. CALVIN: A benchmark  
536 for language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE*  
537 *Robotics and Automation Letters (RA-L)*, 7(3):7327–7334, 2022.
- 538 [49] Meta AI Research. SAM 3: Segment anything with concepts. *arXiv preprint arXiv:2511.16719*,  
539 2025.
- 540 [50] Microsoft VITRA Team. VITRA: Scalable vision-language-action model pretraining for robotic  
541 manipulation with real-life human activity videos. *arXiv preprint arXiv:2510.21571*, 2025.
- 542 [51] Kaichun Mo, Leonidas Guibas, Mustafa Mukadam, Abhinav Gupta, and Shubham Tulsiani.  
543 Where2act: From pixels to actions for articulated 3d objects. In *International Conference on*  
544 *Computer Vision (ICCV)*, 2021.
- 545 [52] Austin Myers, Ching L. Teo, Cornelia Fermüller, and Yiannis Aloimonos. Affordance detection  
546 of tool parts from geometric features. In *ICRA*, 2015.
- 547 [53] Soroush Nasiriany, Sean Kirmani, Tianli Ding, Laura Smith, Yuke Zhu, Danny Driess, Dorsa  
548 Sadigh, and Ted Xiao. RT-Affordance: Affordances are versatile intermediate representations  
549 for robot manipulation. *CoRR*, abs/2411.02704, 2024. doi: 10.48550/ARXIV.2411.02704. URL  
550 <https://doi.org/10.48550/arXiv.2411.02704>.

- 551 [54] Anh Nguyen, Dimitrios Kanoulas, Darwin G Caldwell, and Nikos G Tsagarakis. Object-based  
552 affordances detection with convolutional neural networks and dense conditional random fields.  
553 In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017.
- 554 [55] Toan Nguyen, Minh Nhat Vu, An Vuong, Dzung Nguyen, Thieu Vo, Ngan Le, and Anh  
555 Nguyen. Open-vocabulary affordance detection in 3D point clouds. In *IEEE/RSJ International  
556 Conference on Intelligent Robots and Systems (IROS)*, 2023.
- 557 [56] Xichen Pan, Satya Narayan Shukla, Aashu Singh, Zhuokai Zhao, Shlok Kumar Mishra, Jialiang  
558 Wang, Zhiyang Xu, Jiu hai Chen, Kunpeng Li, Felix Juefei-Xu, Ji Hou, and Saining Xie. Transfer  
559 between modalities with MetaQueries. *arXiv preprint arXiv:2504.06256*, 2025.
- 560 [57] Shengyi Qian, Weifeng Chen, Min Bai, Xiong Zhou, Zhuowen Tu, and Li Erran Li. Affor-  
561 dancellm: Grounding affordance from vision language models. In *IEEE/CVF Conference on  
562 Computer Vision and Pattern Recognition, CVPR 2024 - Workshops, Seattle, WA, USA, June  
563 17-18, 2024*, pages 7587–7597. IEEE, 2024. doi: 10.1109/CVPRW63382.2024.00754. URL  
564 <https://doi.org/10.1109/CVPRW63382.2024.00754>.
- 565 [58] Daniel Seita, Yufei Wang, Sarthak Shetty, Edward Li, Zackory Erickson, and David Held.  
566 ToolFlowNet: Robotic manipulation with tools via predicting tool flow from point clouds. In  
567 *Conference on Robot Learning (CoRL)*, 2022.
- 568 [59] Yawen Shao, Wei Zhai, Yuhang Yang, Hongchen Luo, Yang Cao, and Zheng-Jun Zha. GREAT:  
569 Geometry-intention collaborative inference for open-vocabulary 3D object affordance grounding.  
570 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition  
571 (CVPR)*, pages 17326–17336, 2025.
- 572 [60] Xiaohao Sun, Hanxiao Jiang, Manolis Savva, and Angel X. Chang. OPDMulti: Openable part  
573 detection for multiple objects. In *International Conference on 3D Vision, 3DV 2024, Davos,  
574 Switzerland, March 18-21, 2024*, pages 169–178. IEEE, 2024. doi: 10.1109/3DV62453.2024.  
575 00100. URL <https://doi.org/10.1109/3DV62453.2024.00100>.
- 576 [61] Bin Tan, Changjiang Sun, Xiage Qin, Hanat Adai, Zelin Fu, Tianxiang Zhou, Han Zhang,  
577 Yinghao Xu, Xing Zhu, Yujun Shen, and Nan Xue. Masked depth modeling for spatial  
578 perception. *arXiv preprint arXiv:2601.17895*, 2026.
- 579 [62] Yihe Tang, Wenlong Huang, Yingke Wang, Chengshu Li, Roy Yuan, Ruohan Zhang, Jiajun  
580 Wu, and Li Fei-Fei. UAD: Unsupervised affordance distillation for generalization in robotic  
581 manipulation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2025.
- 582 [63] Yingbo Tang, Lingfeng Zhang, Shuyi Zhang, Yinuo Zhao, and Xiaoshuai Hao. Roboafford:  
583 A dataset and benchmark for enhancing object and spatial affordance learning in robot ma-  
584 nipulation. In *Proceedings of the 33rd ACM International Conference on Multimedia*, pages  
585 12706–12713, 2025.
- 586 [64] Yingbo Tang, Shuaike Zhang, Xiaoshuai Hao, Pengwei Wang, Jianlong Wu, Zhongyuan Wang,  
587 and Shanghang Zhang. Affordgrasp: In-context affordance reasoning for open-vocabulary task-  
588 oriented grasping in clutter. In *IEEE/RSJ International Conference on Intelligent Robots and  
589 Systems, IROS 2025, Hangzhou, China, October 19-25, 2025*, pages 9433–9439. IEEE, 2025.  
590 doi: 10.1109/IROS60139.2025.11245995. URL [https://doi.org/10.1109/IROS60139.  
591 2025.11245995](https://doi.org/10.1109/IROS60139.2025.11245995).
- 592 [65] Qwen Team. Qwen3.5: Accelerating productivity with native multimodal agents, February  
593 2026. URL <https://qwen.ai/blog?id=qwen3.5>.
- 594 [66] Tongxuan Tian, Xuhui Kang, and Yen-Ling Kuo. O3afford: One-shot 3d object-to-object  
595 affordance grounding for generalizable robotic manipulation. 2025.
- 596 [67] Zifu Wan, Yaqi Xie, Ce Zhang, Zhiqiu Lin, Zihan Wang, Simon Stepputtis, Deva Ramanan,  
597 and Katia Sycara. InstructPart: Task-oriented part segmentation with instruction reasoning.  
598 In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics  
599 (ACL)*, 2025.

- 600 [68] Hanqing Wang, Mingyu Liu, Xiaoyu Chen, Chengwei Ma, Yiming Zhong, Wenti Yin, Yuhao  
601 Liu, Zhiqing Cui, Jiahao Yuan, Lu Dai, Zhiyuan Ma, and Hui Xiong. Videoafford: Grounding  
602 3d affordance from human-object-interaction videos via multimodal large language model,  
603 2026. URL <https://arxiv.org/abs/2602.09638>.
- 604 [69] Hanqing Wang, Shaoyang Wang, Yiming Zhong, Zemin Yang, Jiamin Wang, Zhiqing Cui,  
605 Jiahao Yuan, Yifan Han, Mingyu Liu, and Yuexin Ma. Affordance-R1: Reinforcement learning  
606 for generalizable affordance reasoning in multimodal large language models. In *Proceedings of  
607 the AAAI Conference on Artificial Intelligence (AAAI)*, 2026.
- 608 [70] Xinyi Wang, Xun Yang, Yanlong Xu, Yuchen Wu, Zhen Li, and Na Zhao. AffordBot: 3D  
609 fine-grained embodied reasoning via multimodal large language models. In *Advances in Neural  
610 Information Processing Systems (NeurIPS)*, 2025.
- 611 [71] Yi-Lin Wei, Mu Lin, Yuhao Lin, Jian-Jian Jiang, Xiao-Ming Wu, Ling-An Zeng, and Wei-  
612 Shi Zheng. AffordDexGrasp: Open-set language-guided dexterous grasp with generalizable-  
613 instructive affordance. In *Proceedings of the IEEE/CVF International Conference on Computer  
614 Vision (ICCV)*, 2025.
- 615 [72] Dongming Wu, Yanping Fu, Saikhe Huang, Yingfei Liu, Fan Jia, Nian Liu, Feng Dai, Tiancai  
616 Wang, Rao Muhammad Anwer, Fahad Shahbaz Khan, and Jianbing Shen. RAGNet: Large-scale  
617 reasoning-based affordance segmentation benchmark towards general grasping. In *Proceedings  
618 of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2025.
- 619 [73] Kun Wu, Chengkai Hou, Jiaming Liu, Zhengping Che, Xiaozhu Ju, Zhuqin Yang, Meng Li,  
620 Yinuo Zhao, Zhiyuan Xu, Guang Yang, et al. RoboMIND: Benchmark on multi-embodiment  
621 intelligence normative data for robot manipulation. *arXiv preprint arXiv:2412.13877*, 2024.
- 622 [74] Ruihai Wu, Yan Zhao, Kaichun Mo, Zizheng Guo, Yian Wang, Tianhao Wu, Qingnan Fan,  
623 Xuelin Chen, Leonidas Guibas, and Hao Dong. VAT-mart: Learning visual action trajectory  
624 proposals for manipulating 3d ARTiculated objects. In *International Conference on Learning  
625 Representations*, 2022. URL <https://openreview.net/forum?id=iEx3PiooLy>.
- 626 [75] Shijie Wu, Yihang Zhu, Yunao Huang, Kaizhen Zhu, Jiayuan Gu, Jingyi Yu, Ye Shi, and Jingya  
627 Wang. AffordDP: Generalizable diffusion policy with transferable affordance. In *Proceedings  
628 of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 6971–6980, 2025.
- 629 [76] Xiaoyang Wu, Daniel DeTone, Duncan Frost, Tianwei Shen, Chris Xie, Nan Yang, Jakob Engel,  
630 Richard Newcombe, Hengshuang Zhao, and Julian Straub. Sonata: Self-supervised learning of  
631 reliable point representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
632 and Pattern Recognition (CVPR)*, 2025.
- 633 [77] Peiran Xu and Yadong MU. Weakly-supervised affordance grounding guided by part-level  
634 semantic priors. In *The Thirteenth International Conference on Learning Representations*, 2025.  
635 URL <https://openreview.net/forum?id=0823rvTIhs>.
- 636 [78] Rongtao Xu, Jian Zhang, Minghao Guo, Youpeng Wen, Haoting Yang, Min Lin, Jianzheng  
637 Huang, Zhe Li, Kaidong Zhang, Liqiong Wang, Yuxuan Kuang, Meng Cao, Feng Zheng, and  
638 Xiaodan Liang. AO: An affordance-aware hierarchical model for general robotic manipulation.  
639 In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2025.
- 640 [79] Yuhang Yang, Wei Zhai, Hongchen Luo, Yang Cao, Jiebo Luo, and Zheng-Jun Zha. Ground-  
641 ing 3d object affordance from 2d interactions in images. In *Proceedings of the IEEE/CVF  
642 International Conference on Computer Vision (ICCV)*, pages 10905–10915, October 2023.
- 643 [80] Bo-Wen Yin, Jiao-Long Cao, Ming-Ming Cheng, and Qibin Hou. DFormerv2: Geometry  
644 self-attention for RGBD semantic segmentation. In *Proceedings of the Computer Vision and  
645 Pattern Recognition Conference*, pages 19345–19355, 2025.
- 646 [81] Tomoya Yoshida, Shuhei Kurita, Taichi Nishimura, and Shinsuke Mori. Developing vision-  
647 language-action model from egocentric videos. *CoRR*, abs/2509.21986, 2025. doi: 10.48550/  
648 ARXIV.2509.21986. URL <https://doi.org/10.48550/arXiv.2509.21986>.

- 649 [82] Tomoya Yoshida, Shuhei Kurita, Taichi Nishimura, and Shinsuke Mori. Text-driven affordance  
650 learning from egocentric vision. *Adv. Robotics*, 39(16):1041–1052, 2025. doi: 10.1080/  
651 01691864.2025.2535676. URL <https://doi.org/10.1080/01691864.2025.2535676>.
- 652 [83] Tomoya Yoshida, Shuhei Kurita, Taichi Nishimura, and Shinsuke Mori. Generating 6DoF  
653 object manipulation trajectories from action description in egocentric vision. In *Proceedings*  
654 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages  
655 17370–17382, June 2025.
- 656 [84] Chunlin Yu, Hanqing Wang, Ye Shi, Haoyang Luo, Sibe Yang, Jingyi Yu, and Jingya Wang.  
657 SeqAfford: Sequential 3D affordance reasoning via multimodal large language model. In  
658 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
659 pages 1691–1701, 2025.
- 660 [85] Qiaojun Yu, Siyuan Huang, Xibin Yuan, Zhengkai Jiang, Ce Hao, Xin Li, Haonan Chang, Junbo  
661 Wang, Liu Liu, Hongsheng Li, Peng Gao, and Cewu Lu. Uniaff: A unified representation of  
662 affordances for tool usage and articulation with vision-language models. In *IEEE International*  
663 *Conference on Robotics and Automation, ICRA 2025, Atlanta, GA, USA, May 19-23, 2025*,  
664 pages 8980–8987. IEEE, 2025. doi: 10.1109/ICRA55743.2025.11127736. URL <https://doi.org/10.1109/ICRA55743.2025.11127736>.
- 666 [86] Chengbo Yuan, Chuan Wen, Tong Zhang, and Yang Gao. General flow as foundation affordance  
667 for scalable robot learning. In *Proceedings of The 8th Conference on Robot Learning (CoRL)*,  
668 volume 270 of *Proceedings of Machine Learning Research*, pages 1541–1566. PMLR, 2024.  
669 URL <https://proceedings.mlr.press/v270/yuan25a.html>.
- 670 [87] Wentao Yuan, Jiafei Duan, Valts Blukis, Wilbert Pumacay, Ranjay Krishna, Adithyavairavan  
671 Murali, Arsalan Mousavian, and Dieter Fox. RoboPoint: A vision-language model for spatial  
672 affordance prediction in robotics. In *Proceedings of The 8th Conference on Robot Learning*  
673 *(CoRL)*, volume 270 of *Proceedings of Machine Learning Research*, pages 4005–4020. PMLR,  
674 2024. URL <https://proceedings.mlr.press/v270/yuan25c.html>.
- 675 [88] Fan Zhang and Michael Gienger. Affordance-based robot manipulation with flow matching.  
676 *CoRR*, abs/2409.01083, 2024. doi: 10.48550/ARXIV.2409.01083. URL [https://doi.org/](https://doi.org/10.48550/arXiv.2409.01083)  
677 [10.48550/arXiv.2409.01083](https://doi.org/10.48550/arXiv.2409.01083).
- 678 [89] Harry Zhang, Ben Eisner, and David Held. FlowBot++: Learning generalized articulated  
679 objects manipulation via articulation projection. In Jie Tan, Marc Toussaint, and Kourosh  
680 Darvish, editors, *Conference on Robot Learning, CoRL 2023, 6-9 November 2023, Atlanta,*  
681 *GA, USA, Proceedings of Machine Learning Research*, pages 1222–1241. PMLR, 2023. URL  
682 <https://proceedings.mlr.press/v229/zhang23c.html>.
- 683 [90] Zhou Zhao, Jie Gao, and Dongyuan Zheng. Affordance-guided robotic grasping via multimodal  
684 large language model reasoning. *IEEE Trans Autom. Sci. Eng.*, 23:4088–4100, 2026. doi:  
685 10.1109/TASE.2026.3651854. URL <https://doi.org/10.1109/TASE.2026.3651854>.
- 686 [91] Ziyang Zhao, Ke Fan, He-Yang Xu, Ning Qiao, Bo Peng, Wenlong Gao, Dongjiang Li, and Hui  
687 Shen. AnchorDP3: 3D affordance guided sparse diffusion policy for robotic manipulation.  
688 *CoRR*, abs/2506.19269, 2025. doi: 10.48550/ARXIV.2506.19269. URL [https://doi.org/](https://doi.org/10.48550/arXiv.2506.19269)  
689 [10.48550/arXiv.2506.19269](https://doi.org/10.48550/arXiv.2506.19269).
- 690 [92] He Zhu, Quyu Kong, Kechun Xu, Xunlong Xia, Bing Deng, Jieping Ye, Rong Xiong, and Yue  
691 Wang. Grounding 3d object affordance with language instructions, visual observations and  
692 interactions. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*  
693 *2025, Nashville, TN, USA, June 11-15, 2025*, pages 17337–17346. Computer Vision Foundation  
694 / IEEE, 2025. doi: 10.1109/CVPR52734.2025.01616. URL [https://openaccess.](https://openaccess.thecvf.com/content/CVPR2025/html/Zhu_Grounding_3D_Object_Affordance_with_Language_Instructions_Visual_Observations_and_CVPR_2025_paper.html)  
695 [thecvf.com/content/CVPR2025/html/Zhu\\_Grounding\\_3D\\_Object\\_Affordance\\_](https://openaccess.thecvf.com/content/CVPR2025/html/Zhu_Grounding_3D_Object_Affordance_with_Language_Instructions_Visual_Observations_and_CVPR_2025_paper.html)  
696 [with\\_Language\\_Instructions\\_Visual\\_Observations\\_and\\_CVPR\\_2025\\_paper.html](https://openaccess.thecvf.com/content/CVPR2025/html/Zhu_Grounding_3D_Object_Affordance_with_Language_Instructions_Visual_Observations_and_CVPR_2025_paper.html).

## 697 **Technical Appendices and Supplementary Material**

698 This appendix provides supplementary material supporting the claims in the main paper. We organise  
699 it as follows. **App. A** specifies the limitations and failed case analysis of our method, and social  
700 responsibility of our work. **App. ??** details the implementation of AFUN, including model archi-  
701 tecture, trainable parameters, the Bézier spline curve motion parameterisation, training schedule,  
702 losses, and inference settings. **App. B** describes the five-stage data-collection pipeline that produces  
703 the multi-source affordance dataset, the filtering criteria, the unit and frame conventions, and the  
704 resulting dataset composition. **App. ??** specifies more details of our experiment, with additional  
705 qualitative results.

### 706 **A Limitations, Failed Case Analysis and social responsibility**

707 **Limitations and Failure Cases.** Though our model has the capability to adapt to open-world  
708 images, this adaptation is still limited by training data motion. Since the only way it learns motion is  
709 from our motion data. Thus, for objects with unique motions that do not exist in the training data, it  
710 will not perform well, as it has no source for learning how such motions would work. We show two  
711 examples in Figure ???. There are no Sun Visor or spray bottle examples in our training dataset, and  
712 there are no similar objects for it to relate to. As a result, it cannot accurately predict how the motion  
713 would go regarding these tasks.

714 **Failure cases.** Typical failures arise when the task phrase underspecifies the target object, when  
715 multiple similar objects or symmetric parts are present, when the actionable region is only weakly  
716 visible, or when depth estimates are unreliable. In such cases, the model may localize a plausible but  
717 incorrect part, predict a curve with the right motion type but wrong direction, or produce a smooth  
718 trajectory that is geometrically valid but insufficient for robust robot execution.

719 **Social responsibility.** AFUN is intended as an intermediate perception-and-motion representation  
720 for controlled robot manipulation, not as a standalone safety mechanism. Physical deployment should  
721 keep human supervision, collision checking, force limits, and task-specific safety constraints in the  
722 execution stack. When releasing data or models, we will document source datasets, filtering steps,  
723 and known failure modes so downstream users can evaluate whether the assumptions match their  
724 intended environment.

### 725 **B Dataset Pipeline Details**

726 We turn heterogeneous demonstrations from robot teleoperation, human egocentric video, simulation  
727 data, and real-world scans into a unified affordance dataset through a two-phase pipeline. First, a  
728 dataset-specific *preprocess* module handles raw-format differences and exports a shared per-interval  
729 schema. Then, a dataset-agnostic *mainprocess* consumes this schema to generate training data: a  
730 SAM3 task query, a tracked object mask, depth, a 3D object trajectory, and fitted Bézier spline curve  
731 parameters. This design keeps format handling separate from affordance label extraction, so adding a  
732 new data source only requires writing a new preprocess adapter. Per-source statistics are reported in  
733 Table 7.

#### 734 **B.1 Cross-Dataset Preprocess**

735 Our preprocess layer is built around source-specific adapters that share the same processing logic  
736 and infrastructure. Every adapter writes the same interval-level schema: observation/contact frames,  
737 RGB-D, task language, camera calibration, and the video span. What differs is how these fields are  
738 recovered from each raw dataset. Below we summarize the dataset-specific handling.

739 **AgiBot.** AgiBotWorld-Beta provides action-interval annotations, and we use each annotated action  
740 as one interval. We load the paired head RGB/depth frames for the interval and use the calibrated  
741 non-fisheye head camera.

Table 7: Per-source dataset statistics. *Episodes* refer to the top-level recording units in each source; *intervals* are action-segmented sub-clips of an episode; *views* are per-camera observations of an interval; *fitted curves* are views with a successfully fit Bézier spline curve.

Source	Modality	Episodes	Intervals	Views	Fitted curves
agibot	robot teleop	17,124	25,395	25,395	1,493
droid	robot teleop	47,508	66,212	132,424	7,381
rh20t	robot teleop	3,588	4,512	34,163	4,127
robomind	robot teleop	14,695	22,642	48,826	17,198
robomind2	robot teleop	6,738	8,918	28,663	4,965
hoi4d	human egocentric	1,020	2,165	2,165	1,974
vitra	human egocentric	3,024	1,098,944	1,098,944	129,433
calvin	simulation	181	6,649	6,649	4,002
rlbench	simulation	268	276	1,104	69
scenefun3d	real scan	585	7,027	148,603	52,692
<b>Total</b>		<b>94,731</b>	<b>1,242,740</b>	<b>1,526,936</b>	<b>223,334</b>

742 **DROID.** DROID provides ZED stereo recordings together with synchronized robot and camera  
743 metadata. We decode the left RGB-D stream from each external stereo camera. Calibration comes  
744 from the official patched files when available; otherwise we fall back to the original HDF5 metadata.  
745 We attach language from the patched files and discard runs marked as failures.

746 **RH20T.** For RH20T, manipulation spans are not annotated directly. We infer them from the gripper-  
747 command trace: a new interval starts when the gripper begins to close and ends when it opens again.  
748 We smooth the signal before this step to avoid noisy frame-level decisions. The interval end is  
749 extended slightly after release to preserve post-contact motion. We export only static camera views  
750 and remove intervals that are too short to provide a meaningful motion trajectory.

751 **RoboMIND.** RoboMIND-v1 mixes several robot embodiments, so the main preprocessing issue is  
752 to make their camera layouts and annotations consistent. We identify the static camera views for each  
753 embodiment and discard arm-mounted views. RGB-D streams are decoded and depth is normalized  
754 to a common metric scale. Manipulation intervals come from the per-step language annotations  
755 released with RoboMIND, which provide frame boundaries for each step. We also remove non-rigid  
756 tasks such as cloth, towel, folding, and rope.

757 **RoboMIND-v2.** RoboMIND-v2 combines Tien Kung, Franka, UR5, and Ark recordings, and the  
758 main challenge is that their gripper signals and usable camera views are not encoded in the same way.  
759 We first identify the robot family for each episode and keep only the camera views that can be used  
760 reliably in our pipeline. Manipulation intervals are then recovered from the family- specific gripper  
761 signal; Ark requires an additional range check because its recordings use two different gripper-value  
762 encodings. We fall back to the whole episode when no valid gripper interval is found, and we remove  
763 non-rigid tasks such as cloth, folding, and rope.

764 **HOI4D.** HOI4D differs from robot sources because the camera is moving, but the dataset also  
765 provides precise geometric annotations: action event markers, camera extrinsics, 3D object assets,  
766 object masks, and object pose trajectories. Because these annotations already determine both the  
767 temporal span and the 3D object motion, we handle HOI4D through a custom path rather than the  
768 robot adapter interface. We use event boundaries such as *Reachout*, *Grasp*, and *Pickup* to form the  
769 manipulation interval, and recover the motion label directly from the recorded object pose trajectory  
770 instead of running depth-based mask tracking. Since the provided masks are object-level, we still run  
771 SAM to obtain the part-level mask used by our affordance supervision.

772 **VITRA.** VITRA is our main source of human manipulation clips. One main challenge for egocen-  
773 tric human videos is noisy camera motion, and VITRA provides per-frame SLAM camera intrinsics

774 and extrinsics for EPIC-KITCHENS and Ego4D clips. We therefore use these camera poses to make  
775 the trajectories geometrically usable. When predicted depth is used, a tracked mask point is first  
776 back-projected in the camera frame where it is observed; the per-frame SLAM poses then transform  
777 this 3D point through the world frame into the observation-frame camera. Since VITRA stores  
778 annotations separately from the source videos, we first resolve each annotated frame index back to  
779 the corresponding EPIC-KITCHENS or Ego4D frame. We use contact-index files to reject clips  
780 without real hand-object contact, and express all projected quantities in the observation-frame camera  
781 coordinates.

782 **Calvin.** Calvin is a simulated robot manipulation dataset with task language, RGB-D observations,  
783 and robot-state traces. We use only the static scene camera, and discard gripper-mounted and tactile  
784 views because they do not provide a stable view of the scene. Calvin’s language annotations mark  
785 task-level windows in a continuous rollout, but they are not contact markers. We use them as semantic  
786 anchors, then use the gripper-action trace to refine the temporal span: an annotation is kept only when  
787 a nearby gripper-closing segment is found, and the final interval covers both the language window  
788 and the matched interaction motion. We use the static RGB-D observation for 3D back-projection,  
789 but do not use simulator-only object variables, such as object poses, drawer states, or button states, as  
790 supervision.

791 **RLBench.** RLBench is a simulated robot manipulation dataset with task language, RGB-D ob-  
792 servations, robot poses, and a binary gripper-open signal. We keep the static scene cameras and  
793 exclude the wrist camera because it moves with the arm. RLBench provides task language at the  
794 episode level rather than per-step temporal spans. We therefore recover intervals from the gripper  
795 signal when possible: closed-gripper segments become manipulation intervals, with the observation  
796 frame shifted earlier to include the approach. Simulator depth is converted to metric depth for 3D  
797 back-projection, but we do not use privileged simulator outputs, such as ground-truth masks or object  
798 states, to generate affordance labels.

799 **SceneFun3D.** SceneFun3D provides posed ARKit RGB-D views of scanned rooms. Each task is  
800 tied to an annotated 3D affordance region and a motion primitive. rather than to a robot or human  
801 action interval. We therefore bypass the tracking, projection, and curve fitting stages in our data  
802 pipeline for SceneFun3D. We sampled the frames at 10 FPS in the dataset videos. We keep steps 2  
803 and 3 (green block in Fig. 2) to retrieve masks of the object of interest, and keep the frames where  
804 that object’s mask is near the starting point of the motion annotation and the projected object 3D  
805 annotation. We then directly convert the motion trajectory to our Bézier spline curve representation,  
806 taking a radius of 90 degrees for circular motion and a fixed 0.3 m length for linear motion.

807 We use a single set of unit and frame conventions throughout: depth maps in millimetres, 3D positions  
808 in metres, and rigid transforms expressed as  $\mathbf{T}_{\text{base} \rightarrow \text{cam}}$ . Every adapter is unit-tested against this  
809 contract before integration.

## 810 B.2 Mainprocess

811 **Query Generation via vLLM Qwen3.5.** The goal of query generation is to produce the text query  
812 that lets SAM3 [49] segment the affordance mask for each task interval. This mask should cover the  
813 object part, visible in the observation image, where contact should be made to execute the task. Since  
814 SAM3 can be driven by text prompts, we first convert the interval’s high-level instruction  $y$  (e.g.,  
815 “Open the cabinet door.”) and visual evidence  $\{I_{\text{obs}}, I_{\text{contact}}\}$  into a short, open-vocabulary  
816 segmentation phrase  $q$ . The phrase is required to name the *minimal manipulable target part*, such  
817 as a handle, knob, button, latch, or lid edge, rather than restating the action or naming the whole  
818 object when a smaller contact part is visible. The observation frame determines how the part should  
819 be described in the target image, while the contact frame helps disambiguate which part is actually  
820 used. The downstream SAM3 video predictor (Sec. B.2) is text-prompted with  $q$ , so the precision of  
821  $q$  directly affects both the affordance mask and the recovered 3D motion.

822 We use Qwen3.5-35B-A3B-FP8 [65] to generate the SAM3 task query. The model returns a JSON  
823 object with a brief rationale and the final `sam3_prompt`. The full system and user prompts are shown  
824 below.

825 [SYSTEM]

826 You are an expert at converting robotics task language + visual  
827 evidence into a compact open-vocabulary segmentation query for SAM3.  
828

829 Project context - Affordance Understanding:  
830 This data is used to train an affordance prediction model. Given an  
831 RGB image and a task-level instruction (e.g., "Open the cabinet  
832 door"), the model must predict (1) the **affordance region** - the  
833 spatial area on the object where physical contact should occur, and  
834 (2) the **post-contact motion trajectory**. The task description  
835 intentionally specifies **what to do**, NOT **where to touch**; the model  
836 must learn the functional mapping from task intent to contact region.  
837 Therefore, only data samples with clear, physically grounded  
838 manipulation targets and meaningful task-level semantics are valuable  
839 for training.  
840

841 Your goal is to produce a short noun phrase that uniquely identifies  
842 the **smallest manipulable target part** (e.g., "top-left drawer  
843 handle", "right knob", "front latch", "left hinge", "power button")  
844 that the robot will operate in the observation image.  
845 You must be extremely concise and precise. You must not describe  
846 actions; only describe the target object/part to segment. Prefer  
847 concrete part names + discriminative attributes (location, row/column,  
848 color, shape, relative position).

849 [USER]  
850 You are given:

851

- 852 \* 'instruction': the task language instruction for this episode/action  
853 interval
- 854 \* 'obs\_frame': the observation frame image (the first frame of the  
855 action interval)
- 856 \* 'contact\_frame': the contact frame image (the frame at the gripper-  
857 close moment)

858

859 Task:  
860 Generate a segmentation text query 'sam3\_prompt' for SAM3 that will  
861 produce a mask of the **minimal target part** the robot is about to  
862 manipulate **in 'obs\_frame'**.

863

864 Guidelines:

865

- 866 1. Output a **single short noun phrase** (1-10 words) suitable for  
867 open-vocabulary segmentation.
- 868 2. When a smaller part is clearly the interaction target, the phrase  
869 must refer to a **physical part** (handle/knob/lid/button/lever/  
870 edge/latch/hinge/tab/strap/rim) rather than a whole object.
- 871 3. Use 'contact\_frame' to infer the true contact target (where the  
872 gripper touches). Use 'obs\_frame' to phrase it in visible terms.
- 873 4. Preserve and include spatial qualifiers if present or inferable  
874 (e.g., "top-left", "second row right", "front", "leftmost",  
875 "upper", "nearest", "on the right side").
- 876 5. If the instruction is ambiguous, resolve it using visual evidence.  
877 If still ambiguous, choose the most likely minimal part and add  
878 one discriminative qualifier (e.g., color/position).
- 879 6. Avoid verbs and action words (open/pull/push/turn). Avoid pronouns  
880 ("it", "that"). Avoid long descriptions.
- 881 7. Do NOT mention "robot", "gripper", "contact", "frame", "image",  
882 "mask", "SAM", or "segmentation".
- 883 8. If the target is a drawer/door, prefer **handle** or **edge**. If  
884 it's a button/switch, prefer **button**/**switch**. If it's a lid,  
885 prefer **lid tab**/**lid edge**. If it's a black cup, prefer  
886 **black cup handle**.

887

888 Output format (strict):  
889 Return ONLY a JSON object with two keys:

```

890 {"rationale": "<your rationale>", "sam3_prompt": "<your noun phrase>"}
891
892 Where:
893
894 * 'rationale' is a very concise and very brief explanation of your
895 reasoning.
896 * 'sam3_prompt' is the best phrase.
897
898 Hard-Fail Policy (must follow):
899 - You output {"rationale": "<your rationale>", "sam3_prompt": null}
900 ONLY when the instruction and the images are fundamentally
901 incompatible such that selecting a manipulable target part would
902 be guesswork.
903 - "Fundamentally incompatible" means: the instruction refers to an
904 object/affordance category that is not present in obs_frame AND
905 there is no clear gripper-contact target in contact_frame that
906 could plausibly satisfy the instruction.
907 - Do NOT fail just because multiple candidates exist; only fail if
908 it is genuinely impossible to identify any plausible target part.
909
910 Hard-Fail Affordance-Relevance Filter (must follow):
911 - You output {"rationale": "<your rationale>", "sam3_prompt": null}
912 when the task is NOT useful for affordance model training. This
913 includes:
914 1. No clear physical manipulation target: the task does not
915 involve contacting and manipulating a specific, localizable
916 part (e.g., "move to the left", "wait", "look around",
917 "navigate to the kitchen").
918 2. Ambiguous / unresolvable target: even with both images, it
919 is impossible to determine a single, well-defined contact
920 region - e.g., the instruction is too vague ("do something with
921 the stuff on the table") and the images provide no
922 disambiguating evidence.
923 3. Non-rigid / deformable / soft-body object: the target is a
924 deformable, soft, or fabric-like object that lacks a well-
925 defined rigid part structure. Our project focuses on rigid
926 objects with strong part-level affordances (handles, knobs,
927 lids, buttons, etc.). Filter out tasks involving cloth, fabric,
928 towels, rope, dough, sponges, or similar soft bodies (e.g.,
929 "fold the towel", "hang the cloth", "flatten the dough",
930 "squeeze the sponge"). Exception: if the task involves grasping
931 a rigid part OF a soft object (e.g., a zipper pull on a
932 jacket), keep it.
933 4. Trivial / non-functional interaction: the task does not
934 teach a meaningful affordance mapping - e.g., the instruction
935 directly names the exact contact part rather than describing a
936 functional task ("grasp the handle", "touch the knob"), or the
937 task is purely a sensor/state check with no physical
938 manipulation.
939 5. Pick-and-place / whole-object relocation: the task is
940 simply grasping an entire object and moving it to a different
941 location. These do not teach meaningful part-level affordances
942 because (a) the contact region is any graspable surface rather
943 than a specific functional part, and (b) the post-contact
944 trajectory is generic relocation (lift -> translate -> place)
945 rather than a functionally determined motion (pull, rotate,
946 press, flip, slide, etc.). Filter out instructions that
947 describe picking up, moving, transferring, or placing an object
948 from one location to another - e.g., "move the cup to the
949 left", "pick up the apple and put it in the bowl", "place the
950 block on the shelf", "put the bottle on the counter",
951 "transfer A to B", "stack the cubes", "sort the objects into
952 the bin". Exception: keep the task if it requires interacting
953 with a specific functional part to achieve the relocation
954 (e.g., "pick up the pot by its handle" names a functional

```

955       grasp point; "slide the drawer out" involves a handle/edge  
 956       affordance).

957 - You strictly filter. Leave out borderline samples. If a task even  
 958 partially matches one of the above categories, output null. High-  
 959 quality training data is far more valuable than quantity - a noisy  
 960 sample hurts the model more than a missing one. Only output a  
 961 valid sam3\_prompt when you are confident the task has a clear,  
 962 rigid, localizable manipulation target with meaningful task-level  
 963 semantics.

964

965 instruction:  
 966 {instruction}

967 **Curve fit.** For each valid object track, we fit the Bézier supervision used in Sec. 4.1 from the  
 968 recovered 3D mask-centroid trajectory. Let  $\{\mathbf{x}_i\}_{i=1}^N$  denote the back-projected object positions  
 969 ordered by frame. Because these points are affected by depth noise, mask jitter, and occasional  
 970 tracking jumps, we first detect abnormally large temporal steps, down-weight them, and smooth  
 971 local neighbourhoods with a robust geometric median. The smoothed trajectory is then resampled  
 972 into a small set of approximately uniform support points along cumulative arc length, with weights  
 973 inversely proportional to the local spatial spread. We fit a planar constant-curvature primitive to these  
 974 support points by nonlinear least squares with a Cauchy robust loss: the primitive is parameterized by  
 975 an origin  $\mathbf{p}_0$ , an orthonormal frame  $(\mathbf{e}_1, \mathbf{e}_2, \mathbf{n})$ , curvature  $\kappa$ , and monotone arc-length coordinates  $s_i$ ,  
 976 so that  $\hat{\mathbf{x}}(s) = \mathbf{p}_0 + s \operatorname{sinc}(\kappa s)\mathbf{e}_1 + s \operatorname{cosec}(\kappa s)\mathbf{e}_2$ . If  $|\kappa|$  times the fitted arc length is below a small  
 977 threshold, the primitive is snapped to a straight line, which avoids overfitting nearly linear motions.  
 978 The fitted curve is sampled densely and converted to the canonical cubic Bézier target by solving a  
 979 least-squares problem for the two interior control points while fixing the start and end points. We  
 980 store the resulting control points relative to the contact anchor  $\mathbf{P}_0$ , matching the model output in  
 981 Eq. 1; the raw trajectory is preserved only for diagnostics and visualization.

### 982 B.3 Gallery on the AFUN dataset

983 We sampled 48 entries from ourmodeltestset and visualize each as the observation frame with the  
 984 SAM3 affordance mask (red mask + yellow bounding box) and the Bézier spline curve fitted 3D  
 985 trajectory (green curve). Cards are split across Fig. 8 and Fig. 9; the language instruction for each  
 986 sample is shown directly below its image.

### 987 B.4 Converting HOVA-500K for Segmentation Training

988 As described in *Stage 2: End-to-End Training for Affordance Segmentation* of Sec. 4.2, our affor-  
 989 dance segmentation model is trained on HOVA-500K [46], RAGNet [72], InstructPart [67], and  
 990 ReasonAFF [69]. RAGNet, InstructPart, and ReasonAFF already match our required format: an  
 991 RGB image, a task query, and a binary affordance mask. HOVA-500K instead provides point-level  
 992 contact supervision. In our loader, each HOVA sample is normalized to an image, an object noun  
 993 field, a verb field, and a contact point. The verb field is source-dependent: it is empty for 3DOI and  
 994 HANDAL, parsed from the Ego4D action field, and read directly from the EPIC-100 verb field. The  
 995 contact point is recovered from the peak of the Gaussian contact heatmap for 3DOI, Ego4D, and  
 996 HANDAL, and from the mean of annotated contact points for EPIC-100. This is useful affordance  
 997 evidence, but it must be converted to dense masks before training our segmentation decoder.

998 To turn the point annotation into a mask, we run a single-frame version of our object-centric annotation  
 999 pipeline. Qwen3-VL receives the image, the normalized noun/verb fields, and the recovered contact  
 1000 point, and produces a compact part-level prompt for SAM3, such as “drawer handle” or “cup rim”.  
 1001 SAM3 then segments the image with this prompt. We keep a mask only if it is both confident and  
 1002 spatially consistent with the HOVA contact point: among masks with confidence above 0.5, we  
 1003 choose the highest-ranked mask whose centroid is within a 2D Euclidean distance of 7% of the image  
 1004 width from the contact point. Samples without such a mask are rejected.

1005 The part-level SAM3 prompt is used only to obtain the mask; it is not used as the training query, since  
 1006 directly naming the contacted part would leak the answer. We therefore run a second Qwen3-VL pass  
 1007 on the original image and the selected-mask overlay. This pass rewrites the sample into a natural  
 1008 task-level instruction that implies the highlighted region without naming it explicitly. We keep a

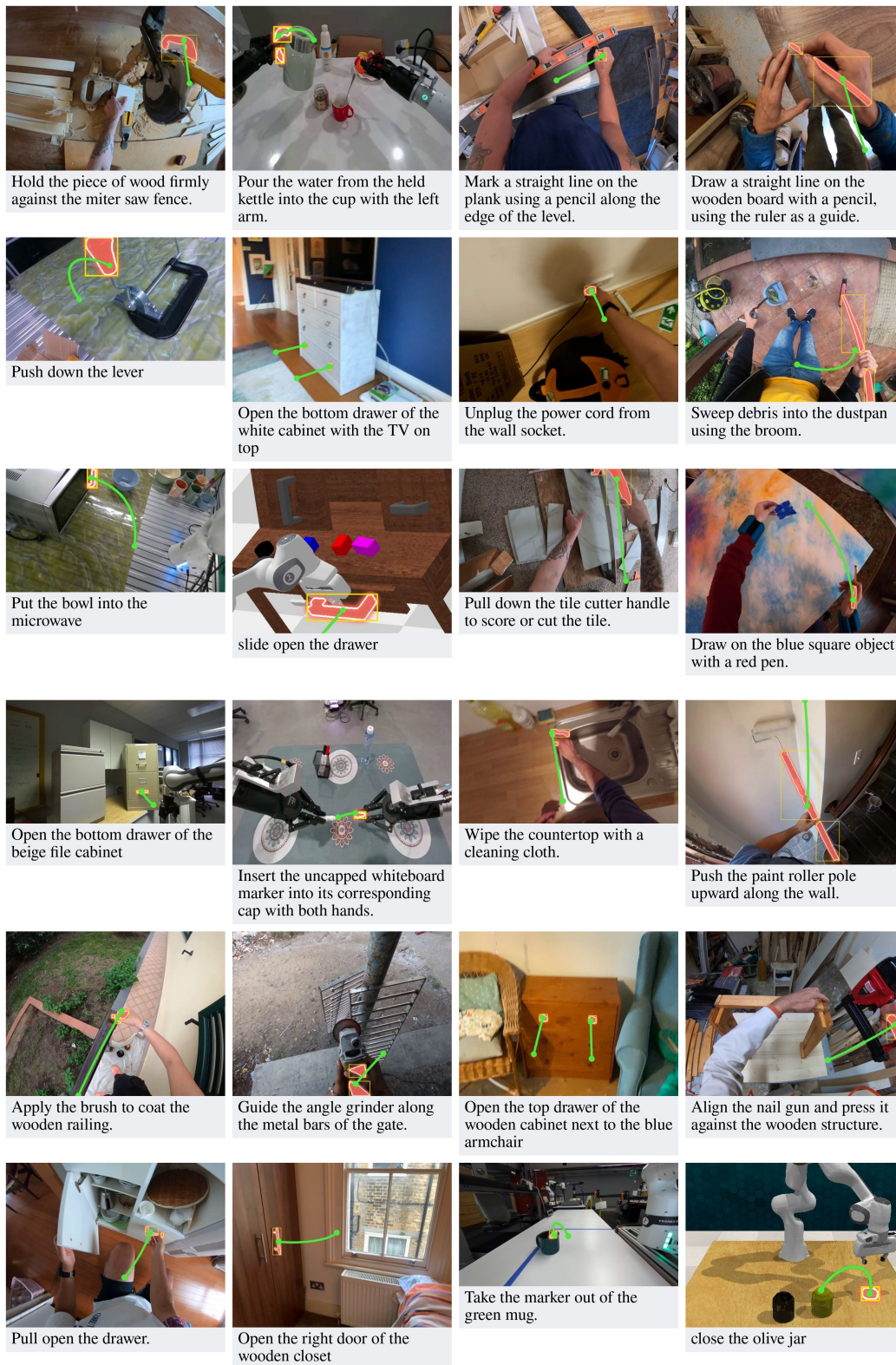


Figure 8: Qualitative gallery on AFUN dataset, Part 1.

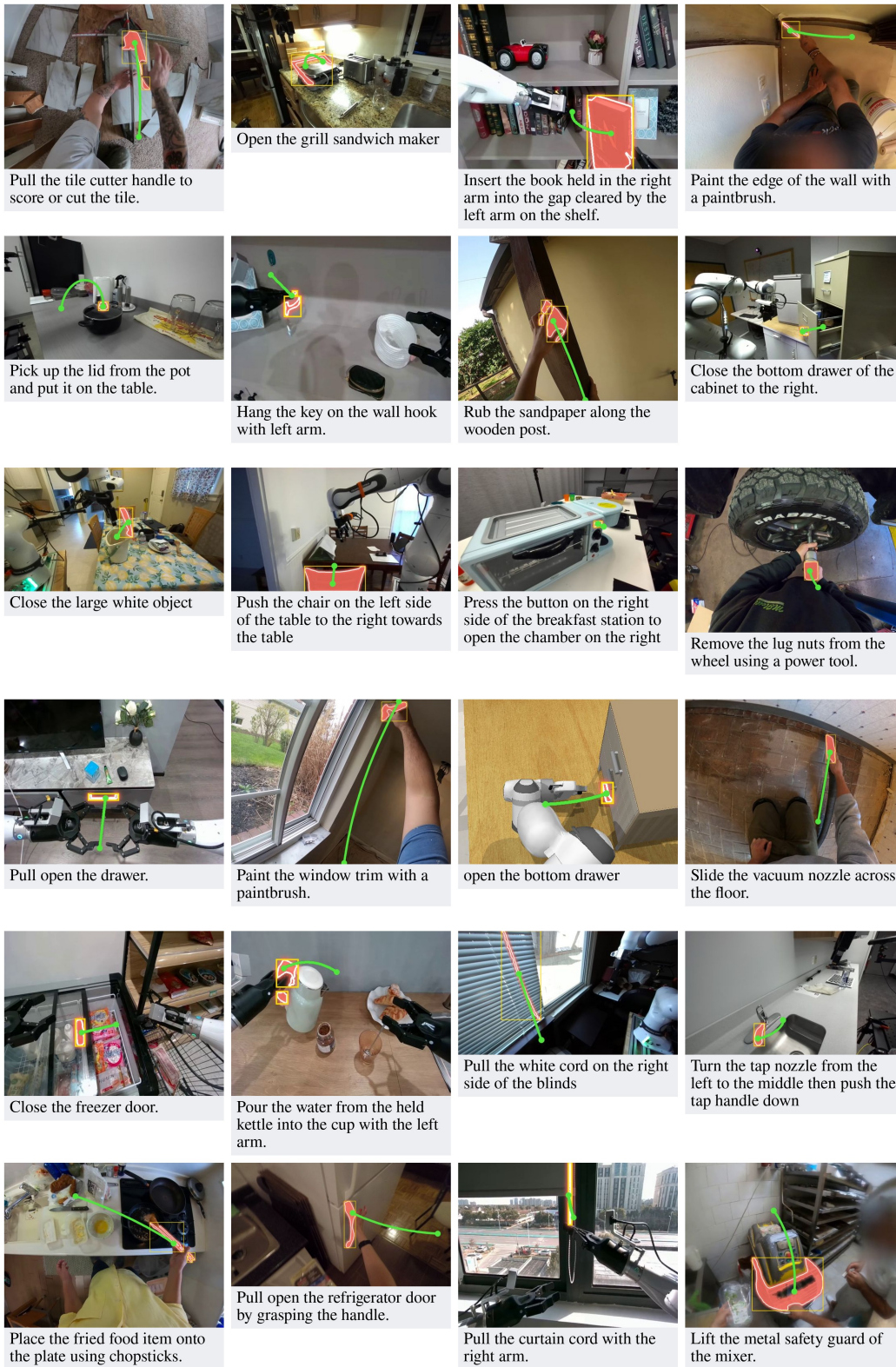


Figure 9: Qualitative gallery on AFUN dataset, Part 2.

1009 sample only when both the mask selection and the rewritten query are valid, resulting in 87,646  
1010 HOVA samples for segmentation training. Random visualization sheets are then inspected for human  
1011 quality verification.

## 1012 C Extended Evaluation Qualitative Results

1013 **Extended Qualitative gallery on the affordance segmentation evaluation.** From Fig. 10 to 13,  
1014 we provide qualitative galleries on the affordance segmentation evaluation results across AFUN,  
1015 AffordanceNet, Affordance-R1, and Qwen3+SAM3. Each row shows the input query, predictions  
1016 from three baseline models, our AFUN, and the ground-truth mask. Red overlays indicate predicted  
1017 regions, green overlays indicate ground truth, and yellow boxes mark mask boxes.

1018 **Qualitative gallery on the 3D motion evaluation.** From Fig. 14 to 18, we additionally provide a  
1019 qualitative gallery for the 3D motion evaluation across AFUN, VRB, VidBot, A0, and General-Flow.  
1020 Each row shows the input task query in the leftmost column, followed by predictions from four  
1021 baseline models, our AFUN, and the ground-truth annotation in the rightmost column. Within each  
1022 cell, the top tile is the predicted trajectory and mask projected onto the input frame, and the bottom  
1023 tile is the back-projected 3D point cloud with the same overlays. Predicted trajectories are coloured  
1024 yellow to blue from start to end; the ground-truth trajectory is rendered as a green curve.

## 1025 D More Qualitative Results

1026 **Extended gallery.** [(zhaoning) 50+ qualitative samples organised by source domain (DROID,  
1027 RH20T, Calvin, Xperience-10M, in-the-wild). Each panel: RGB input + GT mask (green) + AFUN  
1028 mask (blue) + Bézier spline curve curve (orange, blue start / red end).]

1029 **Symmetric and articulated objects.** [(zhaoning) cabinet doors, drawer handles, hinges, faucet  
1030 handles — demonstrate the model picks the correct side and motion direction (open vs. close share  
1031 mask but differ in curve direction).]

1032 **Multi-instance disambiguation.** [(zhaoning) scenes with multiple candidate objects (e.g. two  
1033 mugs, several drawers); the language query selects the correct target.]

1034 **Cross-task on the same object.** [(zhaoning) same scene, two task phrases (e.g. “open the drawer”  
1035 vs. “close the drawer”); show identical mask, opposite Bézier spline curve.]

1036 **Failure cases.** [(zhaoning) honest analysis: occlusion, transparent / glossy objects, deformable  
1037 items, ambiguous language, novel categories. Caption explains the failure mode.]

1038 **Real-robot demonstration.** [(zhaoning) Franka + RGB-D frames with overlay; reference the  
1039 supplementary video. Pair each frame with the AnyGrasp execution step.]

Input query	AffordanceNet	Affordance-R1	Qwen3+SAM3	Ours	GT
Move the dustpan.					
Hold the kitchen strainer.					
Flip the food with the spatula					
Pick up the ladle					
Drink from the mug.					
Vacuum the car floor.					
Use the hammer to drive a nail.					
Pull the tape.					
Use the hammer to strike the wood.					
Pick up the fork					
Turn on the water					

Figure 10: Qualitative gallery on the affordance segmentation evaluation, Part 1.

Input query	AffordanceNet	Affordance-R1	Qwen3+SAM3	Ours	GT
If I want to swipe a surface, which part in the picture would be important?					
If I want to sit on the sofa, which part in the picture should I sit?					
If I want to sit on the sofa, which part in the picture should I sit?					
If I want to sit in the bench, which part in the picture should I sit?					
If I want to dig, which part of the shovel in the picture should be used for containing?					
If I want to cut some food with the knife, which part in the picture should I put the food on?					
If I want to sit on the sofa, which part in the picture should I sit?					
If I want to type, which part in the picture can be utilized?					
If I want to sit in the bench, which part in the picture should I sit?					
If I want to sit in the bench, which part in the picture should I sit?					
If I want to open the toilet, which part in the picture should I interact with?					

Figure 11: Qualitative gallery on the affordance segmentation evaluation, Part 2.

Input query	AffordanceNet	Affordance-R1	Qwen3+SAM3	Ours	GT
If I want to sit in the bench, which part in the picture should I sit?					
If I want to sit in the bench, which part in the picture should I sit?					
If I want to sit in the bench, which part in the picture should I sit?					
If I want to use the fork, which part in the picture should I hold?					
If I want to use the scissors, which part in the picture should I put my fingers in?					
If I want to sit in the bench, which part in the picture should I sit?					
Which part of the mouse is used to scroll the webpage on the computer?					
If I want to sit in the chairs, which part in the picture should I sit?					
Which part is used to cover the pot?					
If I want to sit in the chairs, which part in the picture can support my back?					

Figure 12: Qualitative gallery on the affordance segmentation evaluation, Part 3.































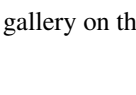
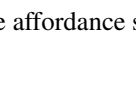
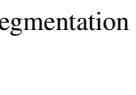
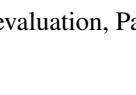
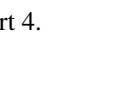
Input query	AffordanceNet	Affordance-R1	Qwen3+SAM3	Ours	GT
To control the cursor on this laptop, where on the device would you place your finger?					
Could you open the left door of the refrigerator? To open the left door of the refrigerator, simply pull on the handle and slide it open.					
Open the left door of the refrigerator. Grasp the handle of the left door and pull it gently to open.					
Fetch me a jug. The handle of the jug is to be grasped for pouring.					
Locate the tap to turn off the water. Grip the tap handle firmly and twist it clockwise to shut off the water flow.					
Can you hand me a knife? The knife can be safely held by its handle for use.					
Locate the tap for me. To operate the tap, twist the handle to control the flow of water.					

Figure 13: Qualitative gallery on the affordance segmentation evaluation, Part 4.

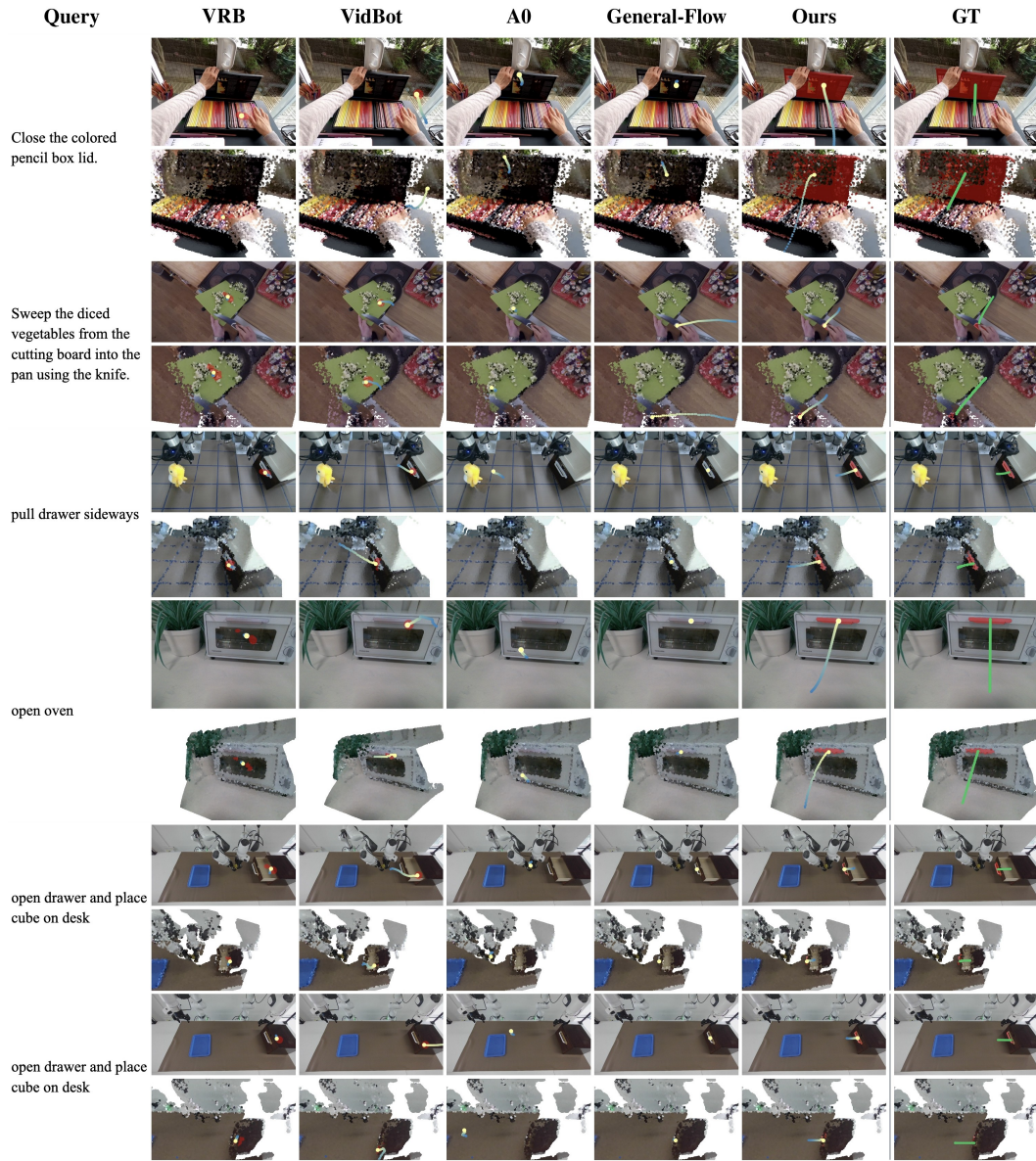


Figure 14: Qualitative gallery on the 3D motion part of the AFUN test set, part 1.

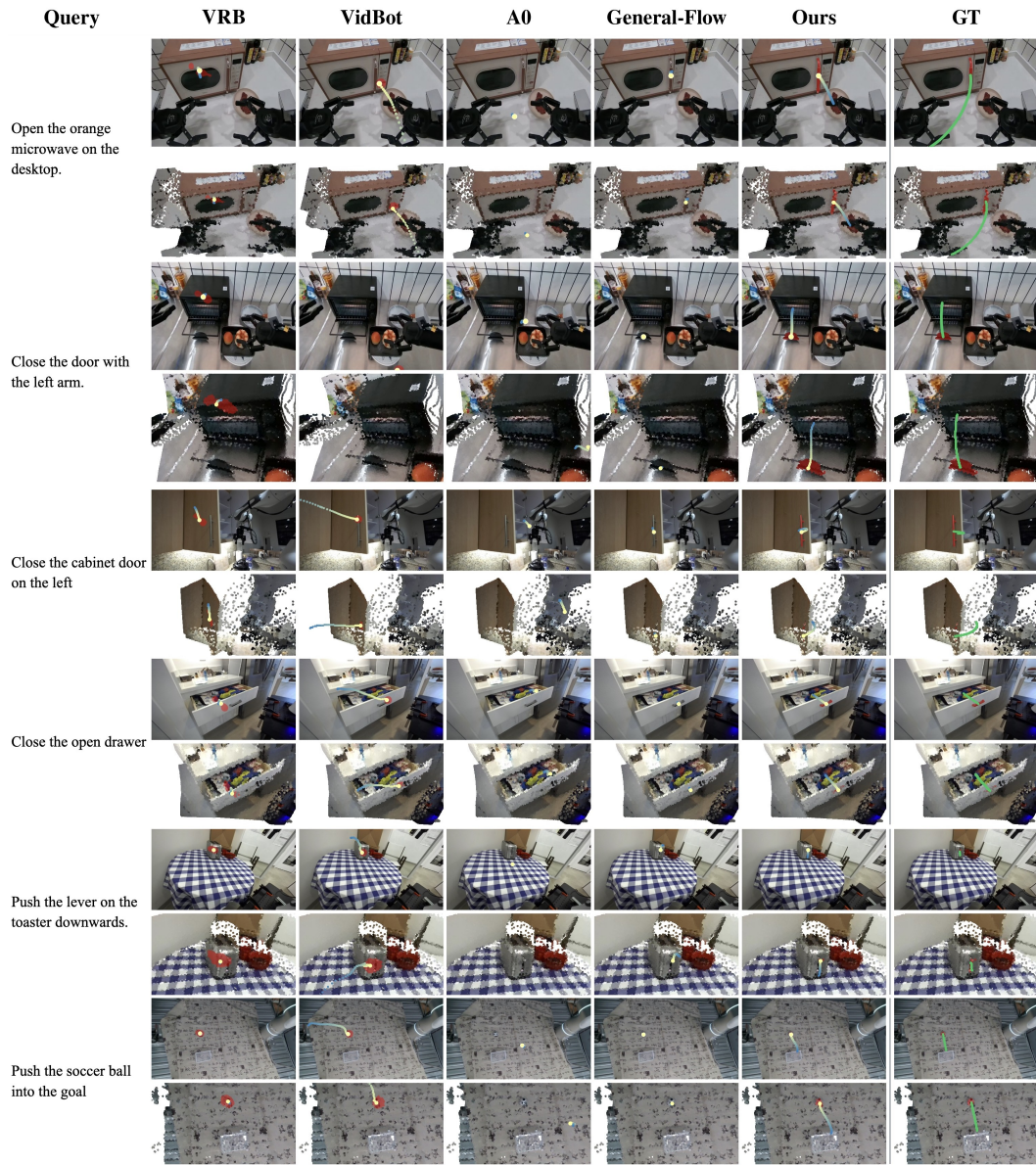


Figure 15: Qualitative gallery on the 3D motion part of the AFUN test set, part 2.

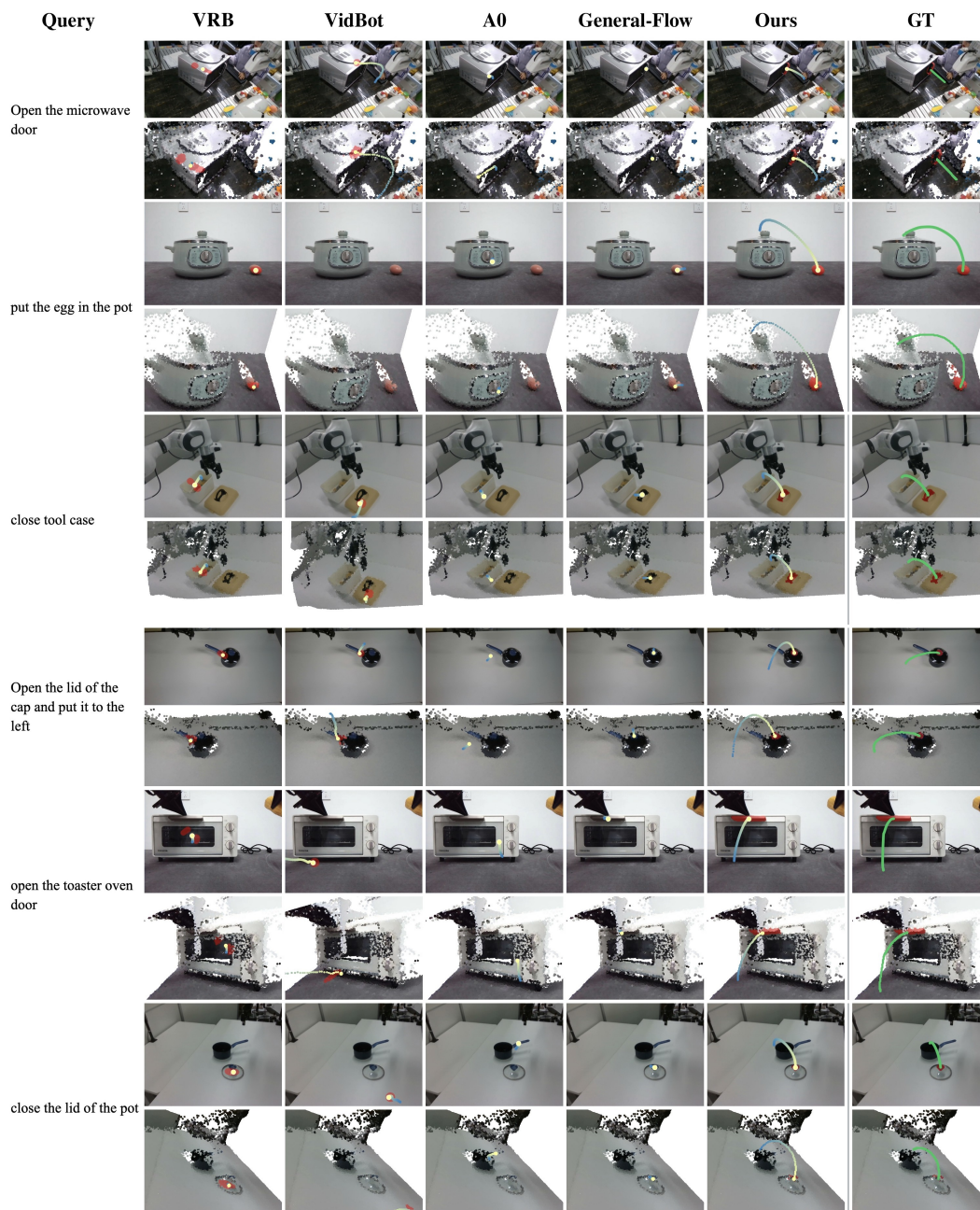


Figure 16: Qualitative gallery on the 3D motion part of the AFUN test set, part 3.

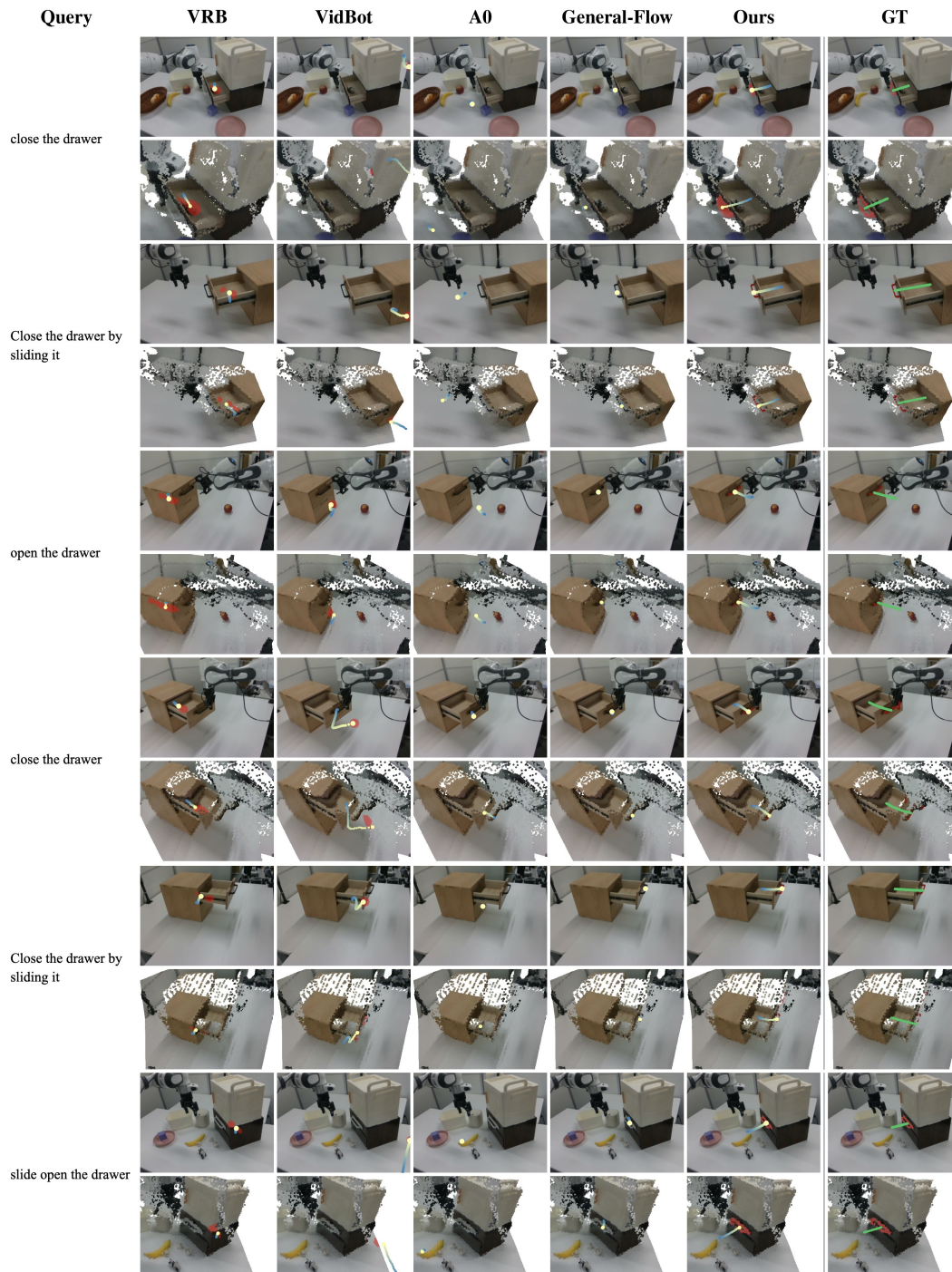


Figure 17: Qualitative gallery on the 3D motion part of the AFUN test set, part 4.

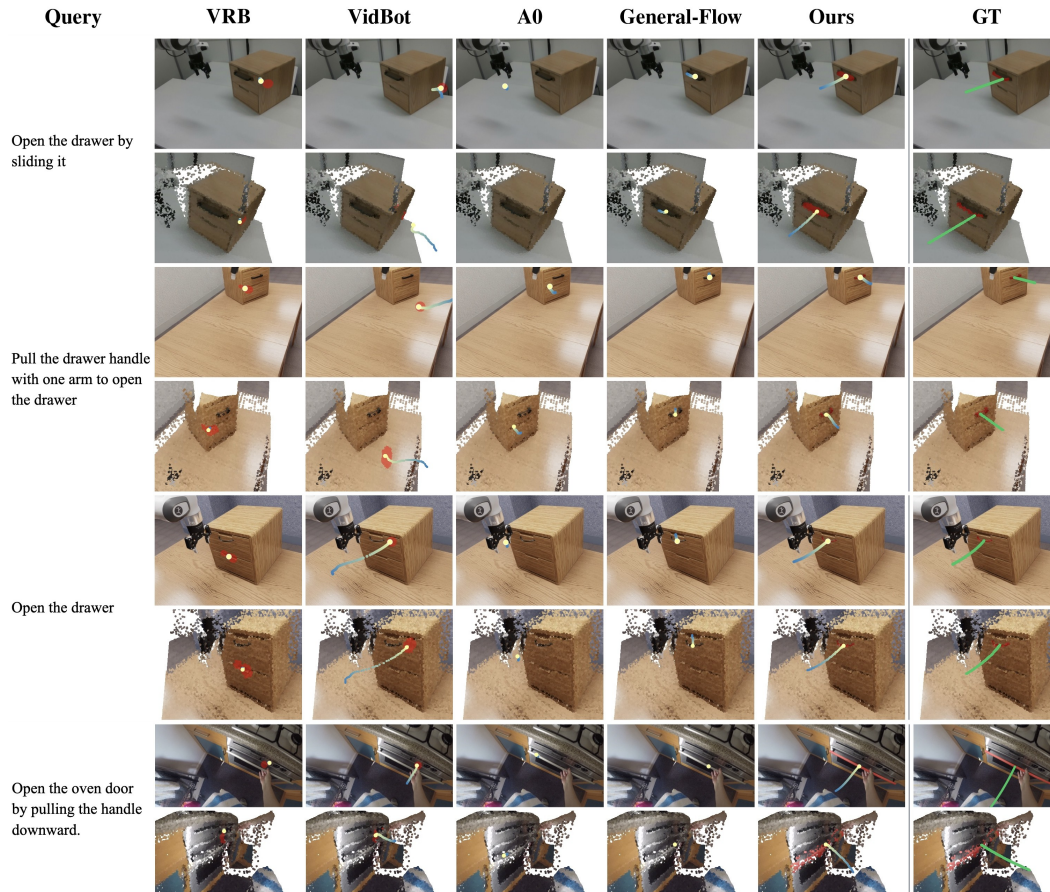


Figure 18: Qualitative gallery on the 3D motion part of the AFUN test set, part 5.

**[Placeholder: Extended Qualitative Gallery]**

6 × 4 grid (or larger). **Rows:** source domains (DROID, RH20T, Calvin, Xperience-10M, in-the-wild, real Franka). **Columns:** (1) RGB input, (2) GT mask (green overlay), (3) AFUN predicted mask (blue overlay), (4) Bézier spline curve (orange) with start (blue) and end (red). Include captions identifying symmetric / articulated / multi-instance cases and mark failure rows in red.

Figure 19: **Extended qualitative results across domains.** AFUN localises the functional region and produces smooth, on-object curves across diverse embodiments and viewpoints, including symmetric parts, articulation, and multi-instance scenes. Failure rows highlight remaining limitations (occlusion, transparency, deformable objects).

1040 **NeurIPS Paper Checklist**

1041 The checklist is designed to encourage best practices for responsible machine learning research,  
1042 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
1043 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
1044 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
1045 towards the page limit.

1046 Please read the checklist guidelines carefully for information on how to answer these questions. For  
1047 each question in the checklist:

- 1048 • You should answer [Yes], [No], or [N/A].
- 1049 • [N/A] means either that the question is Not Applicable for that particular paper or the  
1050 relevant information is Not Available.
- 1051 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

1052 **The checklist answers are an integral part of your paper submission.** They are visible to the  
1053 reviewers, area chairs, senior area chairs, and ethics reviewers. You will also be asked to include it  
1054 (after eventual revisions) with the final version of your paper, and its final version will be published  
1055 with the paper.

1056 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
1057 While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a  
1058 proper justification is given (e.g., error bars are not reported because it would be too computationally  
1059 expensive” or “we were unable to find the license for the dataset we used”). In general, answering  
1060 [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we  
1061 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
1062 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
1063 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
1064 please point to the section(s) where related material for the question can be found.

1065 **IMPORTANT, please:**

- 1066 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 1067 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 1068 • **Do not modify the questions and only use the provided macros for your answers.**

1069 **1. Claims**

1070 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1071 paper’s contributions and scope?

1072 Answer: [Yes]

1073 Justification: We state our main contributions in the abstract and §1, and support them with  
1074 the data, method, and experiments in §3–§5.

1075 Guidelines:

- 1076 • The answer [N/A] means that the abstract and introduction do not include the claims  
1077 made in the paper.
- 1078 • The abstract and/or introduction should clearly state the claims made, including the  
1079 contributions made in the paper and important assumptions and limitations. A [No] or  
1080 [N/A] answer to this question will not be perceived well by the reviewers.
- 1081 • The claims made should match theoretical and experimental results, and reflect how  
1082 much the results can be expected to generalize to other settings.
- 1083 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1084 are not attained by the paper.

1085 **2. Limitations**

1086 Question: Does the paper discuss the limitations of the work performed by the authors?

1087 Answer: [Yes]

1088 Justification: We discuss the main limitations and failure cases in App. A.

1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141

Guidelines:

- The answer [N/A] means that the paper has no limitation while the answer [No] means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate “Limitations” section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren’t acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

**3. Theory assumptions and proofs**

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We do not make new theoretical claims; the formulas in §4 define the representation and training losses we use.

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

**4. Experimental result reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We describe the data pipeline, model, training stages, metrics, baselines, and evaluation protocols in §3–§5 and App. B–??.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.

- 1142
- 1143
- 1144
- 1145
- 1146
- 1147
- 1148
- 1149
- 1150
- 1151
- 1152
- 1153
- 1154
- 1155
- 1156
- 1157
- 1158
- 1159
- 1160
- 1161
- 1162
- 1163
- 1164
- 1165
- 1166
- 1167
- 1168
- 1169
- 1170
- 1171
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
  - If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
  - Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
  - While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
    - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
    - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
    - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
    - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 1172 5. Open access to data and code

1173 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
1174 tions to faithfully reproduce the main experimental results, as described in supplemental  
1175 material?

1176 Answer: [Yes]

1177 Justification: We provide the pipeline, data-source, and evaluation details in the paper and  
1178 will release the code and processed assets with reproduction instructions.

1179 Guidelines:

- 1180
- 1181
- 1182
- 1183
- 1184
- 1185
- 1186
- 1187
- 1188
- 1189
- 1190
- 1191
- 1192
- 1193
- 1194
- 1195
- 1196
- The answer [N/A] means that paper does not include experiments requiring code.
  - Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
  - While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
  - The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
  - The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
  - The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
  - At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

- 1197 • Providing as much information as possible in supplemental material (appended to the  
1198 paper) is recommended, but including URLs to data and code is permitted.

## 1199 6. Experimental setting/details

1200 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
1201 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1202 Answer: [Yes]

1203 Justification: We report the main training setup in §4.2 and §5.1, and give the test splits,  
1204 metrics, and baseline protocols in §5 and App. ??.

1205 Guidelines:

- 1206 • The answer [N/A] means that the paper does not include experiments.
- 1207 • The experimental setting should be presented in the core of the paper to a level of detail  
1208 that is necessary to appreciate the results and make sense of them.
- 1209 • The full details can be provided either with the code, in appendix, or as supplemental  
1210 material.

## 1211 7. Experiment statistical significance

1212 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1213 information about the statistical significance of the experiments?

1214 Answer: [Yes]

1215 Justification: We report multiple task metrics in §5 showing statistically significant improve-  
1216 ments over the baselines, with the evaluation protocol and metric definitions detailed in  
1217 App. ??.

1218 Guidelines:

- 1219 • The answer [N/A] means that the paper does not include experiments.
- 1220 • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
1221 intervals, or statistical significance tests, at least for the experiments that support the  
1222 main claims of the paper.
- 1223 • The factors of variability that the error bars are capturing should be clearly stated (for  
1224 example, train/test split, initialization, random drawing of some parameter, or overall  
1225 run with given experimental conditions).
- 1226 • The method for calculating the error bars should be explained (closed form formula,  
1227 call to a library function, bootstrap, etc.)
- 1228 • The assumptions made should be given (e.g., Normally distributed errors).
- 1229 • It should be clear whether the error bar is the standard deviation or the standard error  
1230 of the mean.
- 1231 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
1232 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
1233 of Normality of errors is not verified.
- 1234 • For asymmetric distributions, the authors should be careful not to show in tables or  
1235 figures symmetric error bars that would yield results that are out of range (e.g., negative  
1236 error rates).
- 1237 • If error bars are reported in tables or plots, the authors should explain in the text how  
1238 they were calculated and reference the corresponding figures or tables in the text.

## 1239 8. Experiments compute resources

1240 Question: For each experiment, does the paper provide sufficient information on the com-  
1241 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
1242 the experiments?

1243 Answer: [Yes]

1244 Justification: We report the main training compute in §5.1, including the 4× NVIDIA  
1245 GH200 setup, training time, stages, and effective batch sizes.

1246 Guidelines:

- 1247 • The answer [N/A] means that the paper does not include experiments.

- 1248
- 1249
- 1250
- 1251
- 1252
- 1253
- 1254
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
  - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 1255 9. Code of ethics

1256 Question: Does the research conducted in the paper conform, in every respect, with the  
1257 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1258 Answer: [Yes]

1259 Justification: We use cited public assets, run robot experiments in controlled settings, and  
1260 discuss social responsibility in App. A.

1261 Guidelines:

- 1262
- 1263
- 1264
- 1265
- 1266
- 1267
- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 1268 10. Broader impacts

1269 Question: Does the paper discuss both potential positive societal impacts and negative  
1270 societal impacts of the work performed?

1271 Answer: [Yes]

1272 Justification: We discuss the potential societal impacts of our work in App. A.

1273 Guidelines:

- 1274
- 1275
- 1276
- 1277
- 1278
- 1279
- 1280
- 1281
- 1282
- 1283
- 1284
- 1285
- 1286
- 1287
- 1288
- 1289
- 1290
- 1291
- 1292
- 1293
- 1294
- 1295
- The answer [N/A] means that there is no societal impact of the work performed.
  - If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 1296 11. Safeguards

1297 Question: Does the paper describe safeguards that have been put in place for responsible  
1298 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
1299 image generators, or scraped datasets)?

1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1350  
1351

Answer: [N/A]

Justification: We do not release a general-purpose language model, image generator, or other high-misuse-risk model; our release is focused on affordance data and models.

Guidelines:

- The answer [N/A] means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cite the datasets, pretrained models, and baselines we use, and we follow their licenses and terms of use.

Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We document the new dataset and model in §3, §4, and App. B, including preprocessing, filtering, curve fitting, training, and evaluation.

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

- 1352 **14. Crowdsourcing and research with human subjects**
- 1353 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1354 include the full text of instructions given to participants and screenshots, if applicable, as  
1355 well as details about compensation (if any)?
- 1356 Answer: [N/A]
- 1357 Justification: We do not conduct new crowdsourcing or human-subject studies; we only use  
1358 existing public human egocentric datasets cited in §3.
- 1359 Guidelines:
- 1360 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1361 with human subjects.
  - 1362 • Including this information in the supplemental material is fine, but if the main contribu-  
1363 tion of the paper involves human subjects, then as much detail as possible should be  
1364 included in the main paper.
  - 1365 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1366 or other labor should be paid at least the minimum wage in the country of the data  
1367 collector.
- 1368 **15. Institutional review board (IRB) approvals or equivalent for research with human**  
1369 **subjects**
- 1370 Question: Does the paper describe potential risks incurred by study participants, whether  
1371 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1372 approvals (or an equivalent approval/review based on the requirements of your country or  
1373 institution) were obtained?
- 1374 Answer: [N/A]
- 1375 Justification: We do not collect new human-subject data, and we rely on the original  
1376 documentation for the public human-video datasets we use.
- 1377 Guidelines:
- 1378 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1379 with human subjects.
  - 1380 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1381 may be required for any human subjects research. If you obtained IRB approval, you  
1382 should clearly state this in the paper.
  - 1383 • We recognize that the procedures for this may vary significantly between institutions  
1384 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1385 guidelines for their institution.
  - 1386 • For initial submissions, do not include any information that would break anonymity (if  
1387 applicable), such as the institution conducting the review.
- 1388 **16. Declaration of LLM usage**
- 1389 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
1390 non-standard component of the core methods in this research? Note that if the LLM is used  
1391 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
1392 scientific rigor, or originality of the research, declaration is not required.
- 1393 Answer: [Yes]
- 1394 Justification: We describe Qwen3-VL as our frozen vision–language backbone in §4, and  
1395 Qwen/vLLM query generation in §3 and App. B.
- 1396 Guidelines:
- 1397 • The answer [N/A] means that the core method development in this research does not  
1398 involve LLMs as any important, original, or non-standard components.
  - 1399 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not  
1400 be described.