

# GLOVER: Generalizable Open-Vocabulary Affordance Reasoning for Task-Oriented Grasping

Teli Ma<sup>1,†</sup>, Zifan Wang<sup>1,†</sup>, Jiaming Zhou<sup>1</sup>, Mengmeng Wang<sup>2</sup>, Junwei Liang<sup>1,3,\*</sup>

<sup>1</sup>AI, HKUST(GZ) <sup>2</sup>ZJUT <sup>3</sup>CSE, HKUST

tma184@connect.hkust-gz.edu.cn junweiliang@hkust-gz.edu.cn

†Equal Contribution \*Corresponding Author

<https://teleema.github.io/projects/GLOVER/>

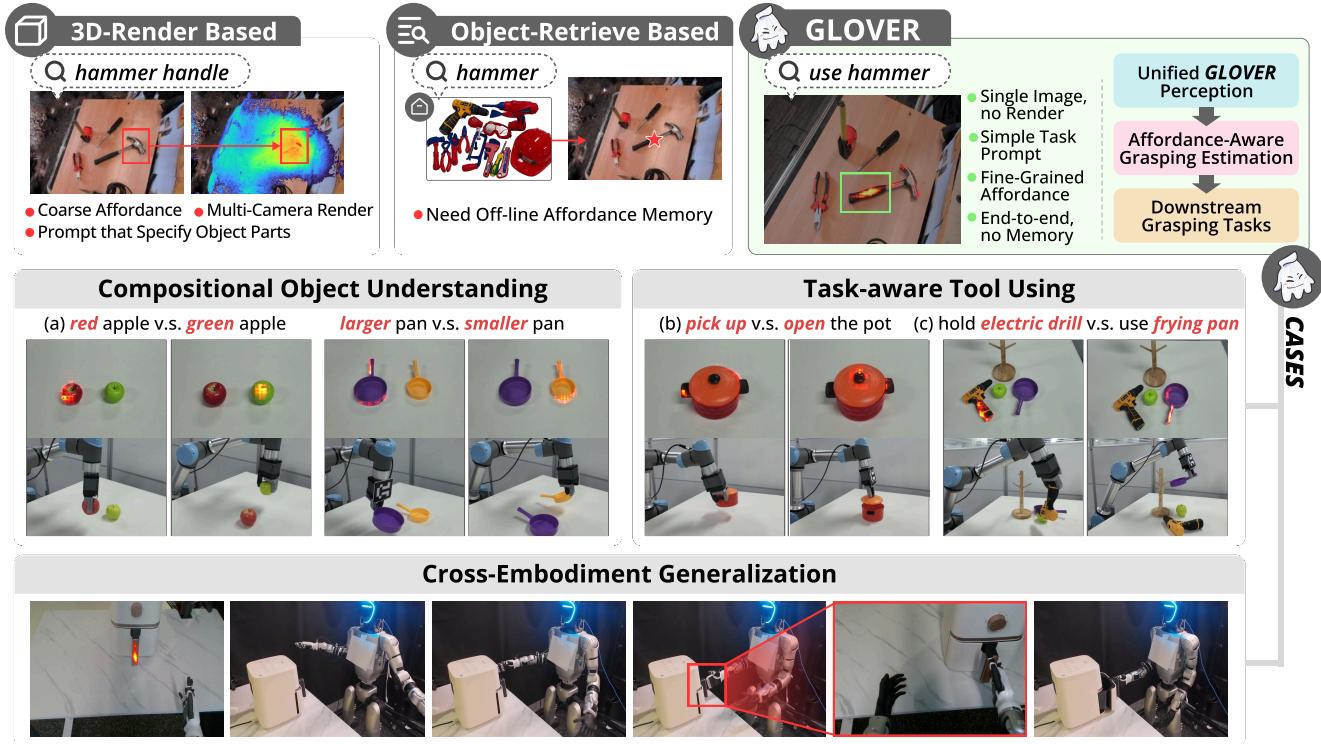


Figure 1. **Top:** Compared with previous methods, GLOVER eliminates the need for capturing multi-view images for 3D rendering, explicit instruction of the grasping part, and building extra off-line affordance memory. GLOVER is capable of providing more accurate fine-grained affordance predictions. **Mid:** We demonstrate the efficacy of GLOVER in handling complex scenarios involving **compositional object understanding** and **task-aware tool using** under targeted human instructions, including: (a) attributes and relations like color, size, material, (b) tool function reasoning according to action, (c) common tool using in complex scenes. **Bottom:** We validate the effectiveness across embodiments (in a humanoid robot with dexterous hands), the robot grasps the oven handle based on inferred affordance.

## Abstract

Inferring affordable (i.e., graspable) parts of arbitrary objects based on human specifications is essential for robots advancing toward open-vocabulary manipulation. Current grasp planners, however, are hindered by limited vision-language comprehension and time-consuming 3D radiance

modeling, restricting real-time, open-vocabulary interactions with objects. To address these limitations, we propose GLOVER, a unified Generalizable Open-Vocabulary Affordance Reasoning framework, which fine-tunes the Large Language Models (LLMs) to predict the visual affordance of graspable object parts within RGB feature space. We compile a dataset of over 10,000 images

*from human-object interactions, annotated with unified visual and linguistic affordance labels, to enable multi-modal fine-tuning. GLOVER inherits world knowledge and common-sense reasoning from LLMs, facilitating more fine-grained object understanding and sophisticated tool-use reasoning. To enable effective real-world deployment, we present Affordance-Aware Grasping Estimation (AGE), a non-parametric grasp planner that aligns the gripper pose with a superquadric surface derived from affordance data. In evaluations across 30 table-top real-world scenes, GLOVER achieves success rates of 86.0% in part identification and 76.3% in grasping, with speeds approximately 29 times faster in affordance reasoning and 40 times faster in grasping pose estimation than the previous state-of-the-art. We also validate the generalization across embodiments, showing effectiveness in humanoid robots with dexterous hands.*

## 1. Introduction

Human beings have an inherent ability to manipulate objects by understanding natural language instructions, such as distinguishing between different object types, identifying objects’ locations, and determining which part to grasp based on the desired task. Motivated by this, research in robotic grasping has evolved from focusing on closed-set objects [16, 52, 71] to open-vocabulary methods [22, 23, 25, 29, 43, 58, 64, 68, 80]. The previous methods mainly consisted of two categories, based on 3D radiance modeling [23, 64, 68, 80] and object retrieving [25, 32] (the top row of Fig. 1). Methods based on 3D radiance modeling necessitate the acquisition of images from multiple cameras and the rendering of each scene, which is time-consuming. On the other hand, approaches based on object retrieving require the establishment of an additional offline object memory, which is labor-intensive. Also, these methods suffer from a lack of complex reasoning capabilities regarding object properties. They face challenges in locating objects with **subtle linguistic distinctions** (e.g., distinguishing between a *red apple* and a *green apple* as shown in Fig. 1) and determining **task-specific graspable parts** (e.g., *pick up the pot* versus *open the pot* as shown in Fig. 1).

In this work, we present a unified Generalizable Open-Vocabulary Affordance Reasoning (GLOVER) framework for open-vocabulary robotic grasping in an end-to-end manner. We aim to leverage the open-vocabulary reasoning capabilities of Large Language Models (LLMs) and fine-tune LLMs to output visual affordance masks. To achieve this, we define the graspable region as a global **visual affordance mask**, inspired by the visual affordance inferring [18, 21, 50, 63] and reasoning segmentation task [33, 44, 67]. Unlike a binary mask, this affordance mask encodes a continuous probability map, representing the likelihood of grasping at various locations. We adopt the af-

fordance mask representation for two reasons: (1) graspable parts are better represented as regions rather than single points, and (2) predicting global masks based on language input aligns more naturally with model decoding [20, 28, 33, 42]. To support this approach, we collect over 10,000 human-object interaction images, using Vision-Language Models (VLMs) to annotate the language labels and unified Gaussian distribution to annotate visual affordance masks.

With the collected dataset, we take the pairs of image and language instruction as input, leveraging LLaVA-7B [42] to perform multi-modal encoding following [33]. The encoded affordance token, which aggregates both visual and linguistic features, is fed into an affordance decoder to output the desired affordance mask. This mask is then projected into 3D space as **stereo affordance** for downstream tasks. The fine-tuning pipeline of GLOVER offers two key benefits: (1) It can leverage extensive 2D human-object interaction images from diverse kinds of datasets [8, 17, 49, 61], overcoming the scarcity of 3D affordance data. (2) The fine-tuning allows GLOVER to inherit world knowledge and common-sense reasoning from the base LLM, enabling compositional object understanding and task-aware tool using in an open-vocabulary manner, as shown in Fig. 1.

To enable real-world grasping, we introduce an Affordance-Aware Grasping Estimation (AGE) module that estimates gripper poses based on the geometry of affordance regions. AGE is a non-parametric method that outperforms learning-based grasping planners [11, 12] in both performance and efficiency (40 times faster), without reliance on additional training data. Inspired by [45, 74], AGE samples grasping poses within the affordance space, determining the target pose by aligning the gripper with the superquadric surface derived from the stereo affordance geometry. The alignment is optimized via nonlinear constrained optimization. This approach eliminates the usage of standalone grasping pose models [25, 32], enabling direct and affordance-aware pose estimation in an end-to-end manner.

To summarize, our contributions are as follows: **(i)** We present GLOVER, a unified end-to-end perception framework designed for open-vocabulary robotic grasping. To enhance its capabilities, we have curated and annotated a dataset of over 10,000 multi-modal affordance images for fine-tuning, enabling GLOVER to leverage world knowledge and common-sense reasoning inherited from large language models (LLMs). **(ii)** We propose an AGE module to perform non-parametric grasping pose estimation, which demonstrates a  $\times 40$  speed improvement over previous approaches. **(iii)** Our GLOVER module achieves state-of-the-art performance in the affordance benchmark, surpassing previous methods by a significant margin. We test the

open-vocabulary grasping capability across 30 challenging table-top real-world scenes, demonstrating an average improvement of **20.0%** in affordance reasoning success rate and **17.3%** in grasping success rate compared to the previous state-of-the-art. We also demonstrate the generalization across diverse scenes and embodiments, showing effectiveness in 4 tasks with humanoid robots and dexterous hands.

## 2. Related Work

### 2.1. Open-Vocabulary Representation for Manipulation

Many recent studies have worked on language-guided robotic tasks like navigation [2, 30, 31] and manipulation [9, 43, 68]. The intervention of language plays a positive role in policy learning [16, 23, 52], value functions estimation [1, 38] and visual perception [64, 80]. Among them, the open-vocabulary manipulation is one of the most significant research topics. Several recent works focus on integrating 2D foundation models with 3D feature fields to achieve a 3D semantic-aware representation for open-vocabulary tasks [64, 68, 70, 80]. These methods distill features from 2D foundation models like CLIP [66], DINO [4] as training objective for NeRF [53] or GaussianSplatting [26] to reconstruct 3D feature fields. In this work, we directly finetune the 2D foundation models to generate the related affordable areas in a generalizable and open-vocabulary manner, which is robust to scenario changes.

### 2.2. Task-Oriented Grasping

Task-oriented grasping refers to grasping different parts of objects based on the tasks. Previous research solves the task via detecting related object parts [7, 22, 36, 48, 56], modeling 3D point clouds for affordance grounding [14, 41, 73], or transferring grasps to new instances based on category [14, 25]. Recent works [9, 22, 25, 37, 48, 68, 73] leverage vision-language models to reason the object parts for grasping. LERF-TOGO [68] derives a rough 3D object mask using DINO [4] features to expand a relevant area locally. Subsequently, a LERF [27] query is conditioned on this mask to separate sub-parts of the object. Robo-ABC [25] reasons the objects’ grasping point by using CLIP [66] to retrieve objects that share semantic similarity from the affordance memory. ShapeGrasp [37] infers the contact points by prompting the large language models via Chain of Thought [76]. Almost all the above methods rely on labeled 3D part-affordance datasets or additional pre-trained grasp networks like GraspNet [11], AnyGrasp [12] to infer grasping pose. However, GLOVER can leverage extensive 2D affordance data, well-trained 2D foundation models, and does not require an additional grasp planner network to estimate poses, which is more efficient.

### 2.3. Visual Affordance Reasoning

Previous research infers the affordance from human-object interactions [19, 21, 50], scene understanding [47, 64, 68] and 3D point cloud grounding [15, 55, 57, 77]. Recently, the foundation models such as LLMs and VLMs, have been integrated to perform affordance reasoning [18, 25, 48, 63, 73, 78]. AffCorrs [18] tackle one-shot visual affordance transfer by querying the object parts to find semantically corresponding ones via pre-trained DINO-ViT [4]. AffordanceLLM [63] train the LLaVA [42] on affordance dataset AGD20K [49], leveraging the world knowledge of the foundation model. Both Robo-ABC [25] and RAM [32] adopt the retrieve-and-transfer framework for zero-shot affordance reasoning. They construct the affordance memory from 2D images and retrieve the similar demonstration from the affordance memory with the help of CLIP [66] to reason affordance in the unseen domain. Our GLOVER does not require the creation of affordance memory, instead, it reasons the affordance leveraging LLM’s world knowledge in an end-to-end manner.

## 3. GLOVER Method

In this section, we tackle two key problems for open-vocabulary affordance reasoning. First, we identify an effective approach to represent visual affordance for finetuning the foundational model (Sec. 3.1). Second, we examine how to integrate affordance knowledge into an existing foundational model while preserving as much of its original world knowledge as possible (Sec. 3.2- 3.4).

### 3.1. VL-Affordance Dataset Construction

**Image Data collection.** To enable affordance reasoning with the world knowledge of LLMs, we leverage the abundant resources of 2D images with visual affordance annotations. We select images from two common affordance datasets, AGD20K [49] and 3DOI [61]. AGD20K is a semi-supervised affordance dataset that includes 23,816 images with 50 categories of objects and 36 categories of affordance, sourced mainly from COCO [40], HICO [5] and free-license websites.

For 3DOI, we use 10,000 images drawn from Articulation [62], EpicKitchen [8], and Taskonomy [79], which contain over 5,000 affordance-related interaction points. However, 3DOI does not include text labels for the interactable objects. Both AGD20K and 3DOI images feature a mix of egocentric and exocentric views. We collect the linguistic and visual affordance annotations as follows.

**Generate the linguistic affordance labels.** To address missing object labels in the 3DOI [61] dataset, we utilize VLMs for linguistic annotating, followed by human cross-validation to correct any errors. Specifically, we first crop objects from the background based on the bounding box an-

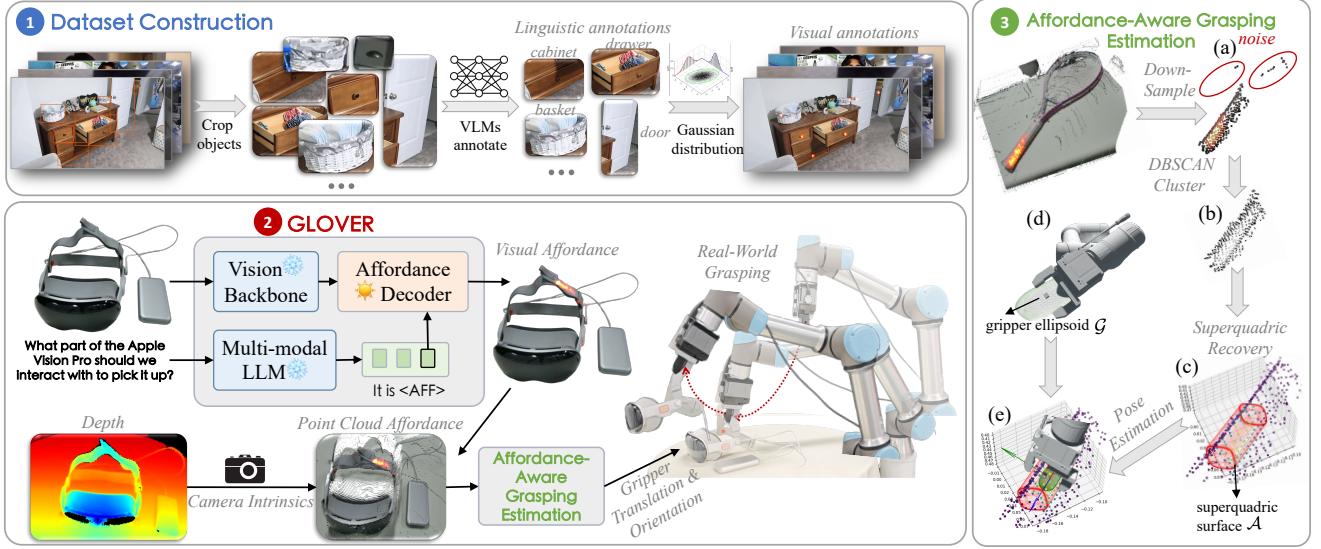


Figure 2. **An overview of our method** (Our contributions are highlighted with colored numbers). 1. We annotate the categories with the VLM and unify the affordance representation. 2. We fine-tune the affordance decoder to decode the affordance token [AFF], which encodes multi-modal information from multi-modal LLM. The fine-tuned GLOVER infers visual affordance in an open-vocabulary manner. 3. The affordance-aware grasping estimation module (AGE), including (a) Voxel down-sampled point clouds. (b) Filter the noise with DBSCAN [69] clustering. (c) Recover superquadric  $\mathcal{A}$  from filtered stereo affordance. (d) Denote the gripper similarly as an ellipsoid surface  $\mathcal{G}$ . (e) Estimate the grasp pose by aligning the  $\mathcal{A}$  and  $\mathcal{G}$ .

notations of the images. The cropped object images are then fed into a vision-language foundation model, SEED-X-17B [13], and prompt it with “*What is the object in the image?*”. We manually filter the generated answers to retain only relevant object categories to ensure accurate annotations. See dataset construction in Fig. 2.

**Unify the visual affordance annotations.** Visual affordances represent the interaction points where humans engage with objects. Predicting affordance through global masks from language inputs aligns more naturally with model decoding. Hence, following the approach in 3DOI [61], we transform affordance points into a 2D Gaussian bump representation. For human-object interaction points in a 2D image, represented as  $a_k = [x_k, y_k]$  for  $k = 1, 2, \dots, N$ , we define the affordance probability mask as:

$$\hat{M}_{aff}(i, j) = \exp\left(-\frac{(i - x_k)^2 + (j - y_k)^2}{2\sigma^2}\right), \quad (1)$$

where  $\sigma$  is the standard deviation, and  $(i, j) \in [0, W] \times [0, H]$  is the spatial indice within the image.

**Dataset overview.** In the end, we obtain a dataset consisting of 12,215 images with 52,240 instances of human-object interactions. The images cover indoor and outdoor scenes, as well as egocentric and exocentric views. Each instance is annotated with both visual and linguistic affordance labels. We believe that our dataset will benefit future research in vision-language affordances.

### 3.2. From Segmentation to Affordance Reasoning

A key distinction between 2D affordance reasoning and traditional 2D instance segmentation is that affordance reasoning produces a continuous probability map, rather than a binary mask. An intuitive idea is to extend VLMs for open-vocabulary segmentation, enabling them to perform open-vocabulary affordance reasoning. This approach enables affordance reasoning while preserving the foundational model’s open-world knowledge in a cost-effective way.

LISA [33] is a large language-instructed segmentation VLM built on an LLM. We adopt the LISA’s structure, which includes a vision backbone, a multi-modal LLM (i.e., LLaVA [42]), and a segmentation decoder. We initialize GLOVER with the LISA-7B pre-trained weights.

### 3.3. Model Details

**Prompt constructing.** To guide the multi-modal LLM (LLaVA [42]) in generating tokens for affordance decoding, we construct prompts in the format: “<IMG> *What part of the [OBJ] should we interact with to [ACT] it?*”, where <IMG> represents image tokens, and [OBJ] and [ACT] specify the object name and action, respectively. “[ACT] it” will be removed from the prompt if the annotation does not exist.

**Multi-modal encoding.** We follow the *Embedding-as-Mask* paradigm in LISA, adding a new affordance token <AFF> to encode combined visual and linguistic features. Given a text prompt  $t$  and an input image  $i$ , we feed them

into the LLaVA model  $\mathcal{F}_{LLM}$  to obtain a response  $r$  (of a sequence of feature vectors):

$$r = \mathcal{F}_{LLM}(i, t), \quad (2)$$

where the encoded <AFF> token feature  $u$  is included in  $r$ . **Visual encoding.** The visual features are important for affordance perception. We adopt the ViT [10] backbone  $\mathcal{F}_{enc}$  to aggregate visual features, encoded as:

$$f = \mathcal{F}_{enc}(i), \quad (3)$$

**Visual affordance decoding.** With the affordance token  $u$  carrying vision-language knowledge, we decode the visual affordance conditioned on it in the visual feature space  $f$ . We follow LISA [33] and SAM [28] to stack Transformer decoder blocks for affordance decoding. Each decoder block consists of self-attention and bi-directional cross-attention. This process is formulated as:

$$M_{aff} = \mathcal{F}_{dec}(u, f), \quad (4)$$

where  $\mathcal{F}_{dec}$  represents the visual affordance decoder. Please refer to Fig. 2 for the pipeline.

### 3.4. Training Objective

To preserve the foundational model’s world knowledge, we freeze the language model component and fine-tune only the affordance decoder parameters specific to our task. Hence, GLOVER is trained end-to-end with only affordance loss to update the affordance decoder. Unlike segmentation mask decoding, which often relies on cross-entropy and DICE losses, we employ the sigmoid focal loss [39] for affordance decoding, as it better handles the continuous distribution of affordance probabilities. The training objective is defined as:

$$\mathbb{L} = \mathcal{L}_{aff} = \text{FL}(M_{aff}, \hat{M}_{aff}). \quad (5)$$

## 4. Affordance-Aware Grasping Estimation

With the inferred visual affordance, the next question is determining how the agent can interact with objects based on these affordances to effectively manipulate them. This question is challenging as the projected stereo affordances (see Sec.4.1) from the visual affordances have complicated and irregular geometric surfaces, making it difficult to identify a global optimal point to grasp. Superquadric recovery [6, 34, 45, 59, 60, 72, 74] is an effective method for estimating superquadric of irregular surfaces, which can be used to estimate the geometry of stereo affordance. Based on the estimated superquadric of affordance region, the grasping pose is calculated via nonlinear constrained optimization. We describe the details below (and illustrated in the third part of Fig. 2).

### 4.1. Stereo Affordance Preprocessing

We first map the deduced visual affordance onto stereo space using the RGB-D camera’s intrinsic parameters. Since the model may infer multiple clusters in the stereo affordance space, we filter out those with low affordance weights. Specifically, we first down-sample the point cloud with voxel down-sampling. Then, we apply DBSCAN [69], a density-based clustering algorithm, to divide the affordances into distinct clusters. We calculate the mean affordance weight for each cluster and discard the 3D point clouds of clusters with low mean weights.

### 4.2. Superquadric Recovery from Stereo Affordance

Superquadrics represent a class of geometric shapes that can model diverse forms. To determine the grasping pose from complex affordance regions, we reconstruct superquadrics from the geometric structures of the affordance point clouds, representing each superquadric as  $\mathcal{A}$ . This process centers on identifying an optimal set of parameters  $\lambda \in \mathbb{R}^{11}$  to maximize alignment between  $N$  stereo affordance points  $a_i = [x_i, y_i, z_i]$  ( $i = 1, \dots, N$ ) and the superquadric surface.

Inspired by superquadric recovery methods [45, 74], we minimize the radial Euclidean distance from the  $a_i$  to the superquadric surface  $\mathcal{A}$  for aligning. To integrate the affordance weight, this process can be formulated as:

$$\min_{\lambda_A} \sum_{i=1}^N \left( \mathcal{W}_i \sqrt{\lambda_V} (F(a_i, \lambda_A) - 1) \right)^2 + \beta V(\lambda_A). \quad (6)$$

The affordance weight  $\mathcal{W}_i$  ensures the inclination towards affordance points with high probability when estimating surface geometries of superquadrics. The inside-outside function  $(F(a_i, \lambda_A) - 1)^2$  aims to minimize the radial Euclidean distance, while the term  $\lambda_V$  is the superquadric volume coefficient, deprecating the expansion of the superquadric volume. To further control the range of the recovered superquadric, we construct a penalty term based on the estimated volume of superquadrics, termed as  $\beta V(\lambda_A)$  to ensure the robustness of the method for noisy point clouds.

### 4.3. Grasp Pose Estimation

We denote the gripper similarly as an ellipsoid surface following previous works [45, 75]. The gripper’s pose is defined by a 7D vector  $x = [x_g, y_g, z_g, q_g^x, q_g^y, q_g^z, q_g^w]$ , where  $(x_g, y_g, z_g)$  are the coordinates of the gripper’s position and  $(q_g^x, q_g^y, q_g^z, q_g^w)$  are its orientation quaternions. For poses  $x$ , our goal is to identify a pose  $\hat{x}$  that allows the gripper ellipsoid  $\mathcal{G}$  to align with the affordance superquadric  $\mathcal{A}$  while satisfying constraints that ensure  $\hat{x}$  is within the gripper’s reach (the third part in Fig. 2).

This can be formalized as the following nonlinear constrained optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \sum_{i=1}^L ((F(\mathbf{p}_i^x, \lambda_{\mathcal{A}}) - 1))^2, \quad (7)$$

subject to:  
 $C_i(\mathbf{c}_i, \mathbf{p}_1^x, \dots, \mathbf{p}_L^x) > 0.$

The cost function Eq.(7) aims to minimize the distance between the affordance superquadric  $\mathcal{A}$  and  $L$  points  $\mathbf{p}_i^x$ , while  $\mathbf{p}_i^x$  are sampled from the closest half of the gripper ellipsoid  $\mathcal{G}$ . This choice prevents the gripper from penetrating the object by ensuring that only the nearest portion of  $\mathcal{G}$  approaches  $\mathcal{A}$ , thus avoiding potential collisions. Constraint terms  $C_i(\cdot)$  are employed to ensure the generation of safe grasping poses, such as avoiding self-collision and collision with environmental obstacles.

Task	Sub-Tasks	#Num	Objects
Compositional Object Understanding	Attribute	7	apple, mug, electric drill, pan, cup, spoon, fork
	Relation	3	mug, cup, pan
	Complex Scene*	5	tape, goggles, tennis racket, electric drill, fruit
Task-Aware Tool Using	Tool Use	10	screwdriver, knife, scissors, hammer, pliers, electric drill, tape measure, saw, pot, pan
	Function Reason	5	knife, charger, sanitizer bottle, pot, tissue box

Table 1. **The specific details of testing scenes we used in the model evaluation.** GLOVER can conduct open-vocabulary affordance reasoning and grasping on nearly any object encountered in daily life. Balancing experimental requirements and practical constraints, we selected specific objects to construct the experiments, with the aim of standardizing the testing of different models’ capabilities. (\* In the Complex Scene, the distractors are from all objects we own, over 5 items per scene.)

## 5. Experiments

In this section, we establish a comprehensive experimental framework to validate the effectiveness and efficiency of GLOVER. We compare with previous methods in both the real-world grasping (Sec. 5.3) and visual affordance reasoning benchmark (Sec. 5.4). Ablations (Sec. 5.5) are also performed to analyze the model components.

### 5.1. Implementation Details

We fine-tuned our GLOVER based on LISA-7B [33], LLaVA [42], and SAM [28] on eight NVIDIA A6000 GPUs for 5 epochs, requiring approximately 18 hours to complete. The initial learning rate is  $5e - 5$  and we use the AdamW [46] optimizer by default. For the real-world experiments, we collect the RGB-D images with an Orbbec

Femto Bolt, with an image size of  $1280 \times 960$  for UR5e robot arm. The humanoid robot, Unitree G1, is utilized for cross-embodiment generalization with more real downstream experiments. The G1 robot is equipped with Inspire dexterous hands RH56DFX, and we use the original Intel RealSense D435i of the robot to capture RGB-D images with the size of  $640 \times 480$ . The details can be found in supplementary material.

### 5.2. Real-World Task Design

We evaluate GLOVER on a wide range of objects from diverse scenes, focusing on two main challenges: *compositional object understanding*, and *task-aware tool using*. The scenes and objects we used are elaborated in the Table 1.

**Compositional object understanding** refers to the agent’s understanding of the fine-grained object attributes like colors, sizes, materials and relations in a compositional way. Here, the agent must select the correct object based on specific human instructions, even when the instructions are morphologically similar. This task is highly challenging, as it requires the agent to correctly interpret affordances associated with each object under potentially confusing instructions.

We design three sub-tasks for compositional object understanding, namely *attributes*, *relations*, and *complex scenes*. The *attributes* task evaluates the agents’ understanding of colors, sizes, and materials of different objects. The *relations* task challenges the agents’ ability to understand spatial relations, like “above” or “below” mug shown in Fig. 3. For *complex scenes*, we assess the model’s capacity for object-level perception by increasing the difficulty for the models to select objects based on language instructions through complex multi-object scenes.

**Task-aware tool using** requires reasoning about which parts of common tools are relevant to various everyday tasks based on vague functional descriptions. We split the task into two sub-tasks, namely *tool use* and *function reason*.

For *tool use*, we prompt GLOVER with phrases composed of verbs plus nouns, such as “*pick up the hammer*”, to assess the agent’s understanding of the general affordances of everyday tools. In contrast, *function reasoning* tests the agent’s understanding of the functions associated with different parts of a tool. This task requires more precise visual perception and sophisticated reasoning, as different verbs applied to the same tool (e.g., “*pick up*” vs. “*pump*” the sanitizer) imply varying affordances. This setup demands intricate common-sense reasoning from the agent, as shown in Fig. 3.

**Evaluation.** For each scene, we vary the object positions and orientations to test it 10 times, and report the success rate. The success rate includes affordance and grasping success rate, referring to the success rate of affordance reasoning and real-world grasping. We manually defined the

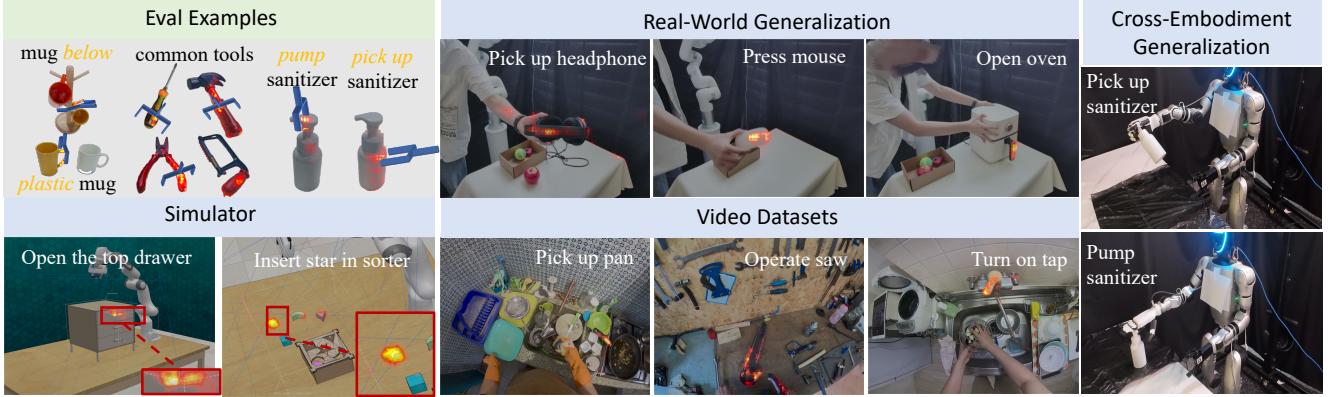


Figure 3. **Eval Examples:** Examples of inferred visual affordance and grasping pose in multiple scenes. The testing scenes are designed to evaluate the model’s compositional object understanding (attributes, relations, complex scenes) and task-aware tools using (tool using, function reasoning). **Generalization:** GLOVER presents open-vocabulary ability across diverse environments, including real-world, simulator (RLBench [24]), scenes from other datasets (Ego4D [17]), and across-embodiments (humanoid robots with dexterous hands).

Method	Compositional Object Understanding (%)						Task-Aware Tool Using (%)				#AVG (%)	
	Attribute(70)		Relation(30)		Complex Scene(50)		Tool Use(100)		Function Reason(50)			
	Aff.	Real	Aff.	Real	Aff.	Real	Aff.	Real	Aff.	Real	Aff.	Real
VRB [3]	52.9	37.1	36.7	20.0	8.0	6.0	58.0	32.0	0.0	0.0	36.7	22.3
LERF-TOGO [68]	58.6	45.7	30.0	23.3	76.0	62.0	64.0	53.0	22.0	10.0	54.7	42.7
LISA* [33]	-	34.0	-	13.3	-	14.0	-	7.0	-	0.0	-	10.7
RAM [32]	64.3	57.1	43.3	26.7	78.0	62.0	84.0	62.0	34.0	12.0	66.0	49.0
<b>GLOVER</b>	<b>95.7</b>	<b>84.3</b>	<b>80.0</b>	<b>73.3</b>	<b>82.0</b>	<b>68.0</b>	<b>88.0</b>	<b>80.0</b>	<b>76.0</b>	<b>68.0</b>	<b>86.0</b>	<b>76.3</b>

Table 2. **Real-world experimental results.** The *Attribute*, *Relation*, *Complex Scene*, *Tool Use* and *Function Reason* have 7, 3, 5, 10, 5 test scenes, respectively. All the scenes are evaluated 10 times by varying the object arrangements. We report the affordance reasoning and real-world grasping success rates to compare with previous methods.

ground truth regions for different objects based on the language instructions of the required task. As long as **the majority of** filtered affordance points fall within these ground truth regions, it is considered an affordance success. Grasping success is defined as the ability to complete real-world grasping. For real-world grasping, we only test the scenes with successfully reasoned affordance.

### 5.3. Real-World Results

**Baselines.** We construct four baselines for comparisons. VRB [3] predicts contact points by learning from human video demonstrations. LERF-TOGO [68] reconstructs the scenes dynamically via LERF [27], extracting 3D object masks from DINO [4] features and conditioning object-part queries based on these masks. We also replace our GLOVER module with the original LISA-7B [33] model, which we denote as LISA\*, to highlight the effectiveness of our affordance fine-tuning. RAM [32] constructs the affordance memory from 2D images and retrieves similar demonstrations from the affordance memory with the help of CLIP [66] to reason affordance in the unseen domain. For a fair comparison, we provide LERF-TOGO with ambiguous language queries rather than specific part queries

(which require an additional LLM to infer). Since LISA\* outputs binary masks, we estimate grasping poses using superquadrics based on the stereo binary mask and report the real-world grasping results.

**Results.** The performance are reported in Table 2. Our model surpasses the previous approach, RAM [32], achieving an average increase of 20.0% in affordance success rate and 17.3% in real-world grasping success rate. Key findings include: (1) VRB [3] performs poorly in scenes that need reasoning based on human instructions due to the lack of language processing capabilities; (2) LERF-TOGO performs well with object recognition but lacks the complex reasoning required for interpreting object relations and diverse tool functions, likely due to the bag-of-words [65] limitation; (3) In contrast to the original LISA model, our model exhibits finer object component perception, resulting in more precise interactive regions and improved real-world grasping accuracy. In summary, GLOVER excels in both intricate common-sense reasoning and precise affordance inference.

## 5.4. Affordance Comparisons

We follow previous SOTA methods to compare affordance capabilities in the following benchmark. AGD20K [49] is a large-scale affordance dataset with a test split for fair comparisons. We follow AffordanceLLM [63] and LOCATE [35] to report the results evaluated on the hard split of AGD20K testing.

**Evaluation metrics.** Following the previous work, we adopt the Kullback-Leibler Divergence (KLD), Similarity (SIM), and Normalized Scanpath Saliency (NSS) as evaluation metrics. Lower KLD values and higher SIM and NSS values indicate better affordance inference. Details on these metrics are elaborated in the supplementary material.

Methods	KLD ↓	SIM ↑	NSS ↑
Cross-View-AG [50]	2.092	0.209	0.138
Cross-View-AG+ [51]	2.034	0.218	0.342
LOCATE [35]	1.829	0.282	0.276
LOCATE-Sup [35]	2.003	0.224	0.435
LOCATE-Sup-OWL [35, 54]	2.127	0.206	0.314
3DOI [61]	4.017	0.200	0.549
VRB [3]	2.154	0.258	0.236
AffordanceLLM [63]	1.661	0.361	0.947
<b>GLOVER</b>	<b>1.098</b>	<b>0.476</b>	<b>1.552</b>

Table 3. **Visual affordance results in the benchmarking dataset.** GLOVER outperforms previous state-of-the-art methods by a large margin in all three metrics.

**Results.** The results in Table 3 show that GLOVER significantly outperforms previous methods across all three metrics. Note that the testing images are **filtered out** from the training set of GLOVER in this comparison. The results highlight the effectiveness of our fine-tuning method in enhancing foundation models for open-vocabulary affordance reasoning.

## 5.5. Ablations

**Efficiency analysis.** We assess GLOVER’s efficiency in terms of time required for affordance reasoning and grasping pose estimation. The process of affordance reasoning includes both scene capture and model inference. As shown in Table 4, GLOVER achieves approximately 330 times and 29 times faster affordance reasoning compared to LERF-TOGO [68] and RAM [32]. For grasping pose estimation, our proposed AGE (Sec. 4) is about 40 times faster than the GraspNet [11] used in LERF-TOGO and RAM. All the time costs are reported on a single NVIDIA 4090 GPU. The results demonstrate the efficiency of our approach. This high efficiency enables GLOVER to track moving objects dynamically, deducing affordances in a real-time way. The tracking performance can be found in the supplementary materials.

**Grasping pose estimation module.** We ablate the grasping pose estimation module to show the effectiveness of the proposed AGE. We evaluate the performance of its modules

Methods	Affordance Inference Time (s)	Grasping Pose Time (s)
LERF-TOGO [68]	~230.0	~4.0
RAM [32]	~20.0	~4.0
<b>GLOVER</b>	<b>~0.7</b>	<b>~0.1</b>

Table 4. The time costs of affordance reasoning and grasping pose process.

from two aspects. The first one is the real-world grasping success rate. Secondly, we assess the pixel-wise spatial distance (PWS-Distance) between the predicted grasp point and the point of maximum probability in the visual affordance map, as a direct measurement of grasping pose quality. We compare AGE with the popular GraspNet [11] and AnyGrasp [12], both of which are learning-based pose estimation models trained on large datasets. Each scene in Sec. 5.2 is tested five times, and we report the average real-world success rate and PWS-Distance across all scenes. Results in Table 5 show that AGE outperforms both learning-based methods in both metrics.

Methods	Real-World SR. (%) ↑	PWS-Distance ↓
GraspNet [11]	62.0	0.179
AnyGrasp [12]	73.3	0.183
<b>AGE</b>	<b>78.7</b>	<b>0.084</b>

Table 5. **The performance of ablating different grasping pose estimation methods.** Real-World SR. and PWS-Distance represent real-world grasping success rate (%) and pixel-wise spatial distance (between [0, 1]), respectively.

**Generalization.** We show the visualization of open-vocabulary affordance reasoning in diverse scenarios to show the generalization in Fig. 3, including *real-world, simulator, other video datasets*. We also validate GLOVER’s generalization ability across embodiments, showing the effectiveness of the proposed method in humanoid robots with dexterous hands. The Inverse Kinematics is utilized to determine the motions for reaching the target grasping pose output by GLOVER. We conduct experiments on four tasks, *open oven, pump sanitizer, pick up sanitizer, and pick up mug*. The results are shown in Table 6. More details and failure cases are elaborated in the supplementary material.

Task	Open Oven	Pump Sanitizer	Pick up Sanitizer	Pick up Mug
<b>GLOVER</b>	4/5	3/5	2/5	3/5

Table 6. Success rate of GLOVER in humanoid robot with dexterous hands.

## 6. Conclusion

In conclusion, our GLOVER framework demonstrates substantial advancements in open-vocabulary affordance reasoning and task-oriented grasping, outperforming existing methods in both accuracy and efficiency. By fine-tuning a foundational model with enhanced affordance understanding while preserving world knowledge, GLOVER achieves

state-of-the-art performance across real-world grasping and affordance reasoning tasks. The AGE module enables rapid and precise grasping pose estimation in a non-parametric manner. Extensive experiments and ablation studies confirm GLOVER’s capabilities in complex reasoning, efficient affordance inference, and robust grasping across diverse scenarios and embodiments.

## References

- [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022. 3
- [2] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3674–3683, 2018. 3
- [3] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances from human videos as a versatile representation for robotics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13778–13790, 2023. 7, 8
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 3, 7
- [5] Yu-Wei Chao, Zhan Wang, Yugeng He, Jiaxuan Wang, and Jia Deng. Hico: A benchmark for recognizing human-object interactions in images. In *Proceedings of the IEEE international conference on computer vision*, pages 1017–1025, 2015. 3
- [6] Laurent Chevalier, Fabrice Jaillet, and Atilla Baskurt. Segmentation and superquadric modeling of 3d objects. 2003. 5
- [7] Fu-Jen Chu, Ruinian Xu, and Patricio A Vela. Learning affordance segmentation for real-world robotic manipulation via synthetic images. *IEEE Robotics and Automation Letters*, 4(2):1140–1147, 2019. 3
- [8] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Evangelos Kazakos, Jian Ma, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision*, pages 1–23, 2022. 2, 3
- [9] Norman Di Palo and Edward Johns. Dinobot: Robot manipulation via retrieval and alignment with vision foundation models. *arXiv preprint arXiv:2402.13181*, 2024. 3
- [10] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 5
- [11] Hao-Shu Fang, Chenxi Wang, Minghao Gou, and Cewu Lu. Graspnet-1billion: A large-scale benchmark for general object grasping. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11444–11453, 2020. 2, 3, 8
- [12] Hao-Shu Fang, Chenxi Wang, Hongjie Fang, Minghao Gou, Jirong Liu, Hengxu Yan, Wenhui Liu, Yichen Xie, and Cewu Lu. Anygrasp: Robust and efficient grasp perception in spatial and temporal domains. *IEEE Transactions on Robotics*, 2023. 2, 3, 8
- [13] Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. Seed-x: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024. 4
- [14] Haoran Geng, Helin Xu, Chengyang Zhao, Chao Xu, Li Yi, Siyuan Huang, and He Wang. Gapartnet: Cross-category domain-generalizable object perception and manipulation via generalizable and actionable parts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7081–7091, 2023. 3
- [15] Yiran Geng, Boshi An, Haoran Geng, Yuanpei Chen, Yaodong Yang, and Hao Dong. Rlafford: End-to-end affordance learning for robotic manipulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5880–5886. IEEE, 2023. 3
- [16] Ankit Goyal, Jie Xu, Yijie Guo, Valts Blukis, Yu-Wei Chao, and Dieter Fox. Rvt: Robotic view transformer for 3d object manipulation. In *Conference on Robot Learning*, pages 694–710. PMLR, 2023. 2, 3
- [17] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 2, 7
- [18] Denis Hadjiveliadis, Sicelukwanda Zwane, Lourdes Agapito, Marc Peter Deisenroth, and Dimitrios Kanoulas. One-shot transfer of affordance regions? affcorrs! In *Conference on Robot Learning*, pages 550–560. PMLR, 2023. 2, 3
- [19] Mahmudul Hassan and Anuja Dharmaratne. Attribute based affordance detection from human-object interaction images. In *Image and Video Technology—PSIVT 2015 Workshops: RV 2015, GPID 2013, VG 2015, EO4AS 2015, MCBMIIA 2015, and VSWS 2015, Auckland, New Zealand, November 23–27, 2015. Revised Selected Papers 7*, pages 220–232. Springer, 2016. 3
- [20] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 2
- [21] Zhi Hou, Baosheng Yu, Yu Qiao, Xiaojiang Peng, and Dacheng Tao. Affordance transfer learning for human-object interaction detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 495–504, 2021. 2, 3

- [22] Siyuan Huang, Haonan Chang, Yuhang Liu, Yimeng Zhu, Hao Dong, Peng Gao, Abdeslam Boularias, and Hongsheng Li. A3vilm: Actionable articulation-aware vision language model. *arXiv preprint arXiv:2406.07549*, 2024. 2, 3
- [23] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023. 2, 3
- [24] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020. 7
- [25] Yuanchen Ju, Kaizhe Hu, Guowei Zhang, Gu Zhang, Min-grun Jiang, and Huazhe Xu. Robo-abc: Affordance generalization beyond categories via semantic correspondence for robot manipulation. *arXiv preprint arXiv:2401.07487*, 2024. 2, 3
- [26] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 3
- [27] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023. 3, 7
- [28] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 2, 5, 6
- [29] Mia Kokic, Danica Kragic, and Jeannette Bohg. Learning task-oriented grasping from human activity datasets. *IEEE Robotics and Automation Letters*, 5(2):3352–3359, 2020. 2
- [30] Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*, pages 104–120. Springer, 2020. 3
- [31] Jacob Krantz, Aaron Gokaslan, Dhruv Batra, Stefan Lee, and Oleksandr Maksymets. Waypoint models for instruction-guided navigation in continuous environments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15162–15171, 2021. 3
- [32] Yuxuan Kuang, Junjie Ye, Haoran Geng, Jiageng Mao, Congyue Deng, Leonidas Guibas, He Wang, and Yue Wang. Ram: Retrieval-based affordance transfer for generalizable zero-shot robotic manipulation. *arXiv preprint arXiv:2407.04689*, 2024. 2, 3, 7, 8
- [33] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 2, 4, 5, 6, 7
- [34] Ales Leonardis, Ales Jaklic, and Franc Solina. Superquadrics for segmenting and modeling range data. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(11):1289–1295, 1997. 5
- [35] Gen Li, Varun Jampani, Deqing Sun, and Laura Sevilla-Lara. Locate: Localize and transfer object parts for weakly supervised affordance grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10922–10931, 2023. 8
- [36] Gen Li, Nikolaos Tsagkas, Jifei Song, Ruaridh Mon-Williams, Sethu Vijayakumar, Kun Shao, and Laura Sevilla-Lara. Learning precise affordances from egocentric videos for robotic manipulation. *arXiv preprint arXiv:2408.10123*, 2024. 3
- [37] Samuel Li, Sarthak Bhagat, Joseph Campbell, Yaqi Xie, Woojun Kim, Katia Sycara, and Simon Stepputtis. Shapegrasp: Zero-shot task-oriented grasping with large language models through geometric decomposition. *arXiv preprint arXiv:2403.18062*, 2024. 3
- [38] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE, 2023. 3
- [39] T Lin. Focal loss for dense object detection. *arXiv preprint arXiv:1708.02002*, 2017. 5
- [40] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 3
- [41] Suhan Ling, Yian Wang, Ruihai Wu, Shiguang Wu, Yuzheng Zhuang, Tianyi Xu, Yu Li, Chang Liu, and Hao Dong. Articulated object manipulation with coarse-to-fine affordance for mitigating the effect of point cloud noise. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 10895–10901. IEEE, 2024. 3
- [42] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 2, 3, 4, 6
- [43] Peiqi Liu, Yaswanth Orru, Chris Paxton, Nur Muhammad Shafiullah, and Lerrel Pinto. Ok-robot: What really matters in integrating open-knowledge models for robotics. *arXiv preprint arXiv:2401.12202*, 2024. 2, 3
- [44] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 2
- [45] Weixiao Liu, Yuwei Wu, Sipu Ruan, and Gregory S Chirikjian. Robust and accurate superquadric recovery: A probabilistic approach. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2676–2685, 2022. 2, 5
- [46] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [47] Guanxing Lu, Shiyi Zhang, Ziwei Wang, Changliu Liu, Jiwén Lu, and Yansong Tang. Manigaussian: Dynamic gaus-

- sian splatting for multi-task robotic manipulation. *arXiv preprint arXiv:2403.08321*, 2024. 3
- [48] Yuhao Lu, Yixuan Fan, Beixing Deng, Fangfu Liu, Yali Li, and Shengjin Wang. VI-grasp: a 6-dof interactive grasp policy for language-oriented objects in cluttered indoor scenes. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 976–983. IEEE, 2023. 3
- [49] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. Learning affordance grounding from exocentric images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2252–2261, 2022. 2, 3, 8
- [50] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. Learning affordance grounding from exocentric images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2252–2261, 2022. 2, 3, 8
- [51] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. Grounded affordance from exocentric view. *International Journal of Computer Vision*, 132(6):1945–1969, 2024. 8
- [52] Teli Ma, Jiaming Zhou, Zifan Wang, Ronghe Qiu, and Junwei Liang. Contrastive imitation learning for language-guided multi-task robotic manipulation. *arXiv preprint arXiv:2406.09738*, 2024. 2, 3
- [53] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 3
- [54] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection. In *European Conference on Computer Vision*, pages 728–755. Springer, 2022. 8
- [55] Kaichun Mo, Leonidas J Guibas, Mustafa Mukadam, Abhinav Gupta, and Shubham Tulsiani. Where2act: From pixels to actions for articulated 3d objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6813–6823, 2021. 3
- [56] Austin Myers, Ching L Teo, Cornelia Fermüller, and Yiannis Aloimonos. Affordance detection of tool parts from geometric features. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1374–1381. IEEE, 2015. 3
- [57] Chuanruo Ning, Ruihai Wu, Haoran Lu, Kaichun Mo, and Hao Dong. Where2explore: Few-shot affordance learning for unseen novel categories of articulated objects. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [58] Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023. 2
- [59] Despoina Paschalidou, Ali Osman Ulusoy, and Andreas Geiger. Superquadrics revisited: Learning 3d shape parsing beyond cuboids. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10344–10353, 2019. 5
- [60] Despoina Paschalidou, Luc Van Gool, and Andreas Geiger. Learning unsupervised hierarchical part decomposition of 3d objects from a single rgb image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1060–1070, 2020. 5
- [61] Shengyi Qian and David F Fouhey. Understanding 3d object interaction from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21753–21763, 2023. 2, 3, 4, 8
- [62] Shengyi Qian, Linyi Jin, Chris Rockwell, Siyi Chen, and David F Fouhey. Understanding 3d object articulation in internet videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1599–1609, 2022. 3
- [63] Shengyi Qian, Weifeng Chen, Min Bai, Xiong Zhou, Zhuowen Tu, and Li Erran Li. Affordancellm: Grounding affordance from vision language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7587–7597, 2024. 2, 3, 8
- [64] Ri-Zhao Qiu, Yafei Hu, Ge Yang, Yuchen Song, Yang Fu, Jianglong Ye, Jiteng Mu, Ruihan Yang, Nikolay Atanasov, Sebastian Scherer, et al. Learning generalizable feature fields for mobile manipulation. *arXiv preprint arXiv:2403.07563*, 2024. 2, 3
- [65] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 7
- [66] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3, 7
- [67] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abderrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13009–13018, 2024. 2
- [68] Adam Rashid, Satvik Sharma, Chung Min Kim, Justin Kerr, Lawrence Yunliang Chen, Angjoo Kanazawa, and Ken Goldberg. Language embedded radiance fields for zero-shot task-oriented grasping. In *7th Annual Conference on Robot Learning*, 2023. 2, 3, 7, 8
- [69] Erich Schubert, Jörg Sander, Martin Ester, Hans Peter Kriegel, and Xiaowei Xu. Dbscan revisited, revisited: why and how you should (still) use dbscan. *ACM Transactions on Database Systems (TODS)*, 42(3):1–21, 2017. 4, 5
- [70] William Shen, Ge Yang, Alan Yu, Jansen Wong, Leslie Pack Kaelbling, and Phillip Isola. Distilled feature fields enable few-shot language-guided manipulation. *arXiv preprint arXiv:2308.07931*, 2023. 3

- [71] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *Conference on Robot Learning*, pages 785–799. PMLR, 2023. 2
- [72] Franc Solina and Ruzena Bajcsy. Recovery of parametric models from range images: The case for superquadrics with global deformations. *IEEE transactions on pattern analysis and machine intelligence*, 12(2):131–147, 1990. 5
- [73] Yaoxian Song, Pinglei Sun, Yi Ren, Yu Zheng, and Yue Zhang. Learning 6-dof fine-grained grasp detection based on part affordance grounding. *arXiv preprint arXiv:2301.11564*, 2023. 3
- [74] Giulia Vezzani, Ugo Pattacini, and Lorenzo Natale. A grasping approach based on superquadric models. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1579–1586. IEEE, 2017. 2, 5
- [75] Giulia Vezzani, Ugo Pattacini, Giulia Pasquale, and Lorenzo Natale. Improving superquadric modeling and grasping with prior on object shapes. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6875–6882. IEEE, 2018. 5
- [76] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022. 3
- [77] Ruihai Wu, Kai Cheng, Yan Zhao, Chuanruo Ning, Guanqi Zhan, and Hao Dong. Learning environment-aware affordance for 3d articulated object manipulation under occlusions. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [78] Kechun Xu, Shuqi Zhao, Zhongxiang Zhou, Zizhang Li, Huaijin Pi, Yifeng Zhu, Yue Wang, and Rong Xiong. A joint modeling of vision-language-action for target-oriented grasping in clutter. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11597–11604. IEEE, 2023. 3
- [79] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3712–3722, 2018. 3
- [80] Yuhang Zheng, Xiangyu Chen, Yupeng Zheng, Songen Gu, Runyi Yang, Bu Jin, Pengfei Li, Chengliang Zhong, Zeng-mao Wang, Lina Liu, et al. Gaussiangrasper: 3d language gaussian splatting for open-vocabulary robotic grasping. *arXiv preprint arXiv:2403.09637*, 2024. 2, 3

## A. Evaluation Metrics

We elaborate on the three metrics we used in evaluating the GLOVER in the affordance benchmark.

**Kullback-Leibler Divergence (KLD)** quantifies the distribution variance between the predicted affordance map  $\mathbf{M}_{aff}$  and the ground truth  $\hat{\mathbf{M}}_{aff}$  ( $\mathbf{M}_{aff}, \hat{\mathbf{M}}_{aff} \in \mathbb{R}^{H \times W}$ ). We first calculate the min-max normalization for

each pixel in the  $\mathbf{M}_{aff}$  and  $\hat{\mathbf{M}}_{aff}$ .

$$\hat{\mathbf{M}}_{aff}^i = \hat{\mathbf{M}}_{aff}^i / \sum \hat{\mathbf{M}}_{aff}, \quad (8)$$

$$\mathbf{M}_{aff}^i = \mathbf{M}_{aff}^i / \sum \mathbf{M}_{aff}. \quad (9)$$

Then the KLD is formulated as:

$$KLD(\hat{\mathbf{M}}_{aff} || \mathbf{M}_{aff}) = \sum_i \hat{\mathbf{M}}_{aff}^i \cdot \log\left(\frac{\hat{\mathbf{M}}_{aff}^i}{\mathbf{M}_{aff}^i}\right). \quad (10)$$

**Similiary (SIM)**, also known as histogram intersection, quantifies the overlap between the predicted affordance map  $\mathbf{M}_{aff}$  and the ground truth  $\hat{\mathbf{M}}_{aff}$ .

$$SIM(\mathbf{M}_{aff}, \hat{\mathbf{M}}_{aff}) = \sum_i \min(\mathbf{M}_{aff}^i, \hat{\mathbf{M}}_{aff}^i). \quad (11)$$

**Normalized Scanpath Saliency (NSS)** evaluates the alignment between the  $\mathbf{M}_{aff}$  and the ground truth  $\hat{\mathbf{M}}_{aff}$ . We first pre-process the  $\mathbf{M}_{aff}$  and  $\hat{\mathbf{M}}_{aff}$  as:

$$\hat{\mathcal{M}} = \mathbb{I}(\hat{\mathbf{M}}_{aff} > 0.1), \quad (12)$$

$$\mathcal{M} = \frac{\mathbf{M}_{aff} - \mu(\mathbf{M}_{aff})}{\sigma(\mathbf{M}_{aff})}, \quad (13)$$

where  $\mathbb{I}$  is the indicator function and  $\mu, \sigma$  represent the mean and standard deviation of  $\mathbf{M}_{aff}$ . NSS is calculated as the mean of the normalized predictions at binary ground truth locations:

$$NSS(\mathcal{M}, \hat{\mathcal{M}}) = \frac{1}{\sum \hat{\mathcal{M}}} \sum_i (\mathcal{M} \times \hat{\mathcal{M}}_i). \quad (14)$$

## B. Real-World Experiment Settings

We introduce the settings of our real-world experiments, including the single table-top robotic arm experiments and the humanoid robot experiments.

For the single table-top robotic arm, we utilize a UR5e robotic arm equipped with a DH PGI gripper and set an Orbbec Femto Bolt on the front side of the workspace. The default image size of the captured RGB-D stream is  $1280 \times 960$ . The setting is shown as Fig. 4.

For the humanoid robot, we use the Unitree G1 equipped with Inspire RH56DFX dexterous hands and Intel RealSense D435i RGB-D camera. The image size for the RGB-D images is set as  $640 \times 480$ . The setting is shown as Fig. 5.

More experiments can be found in Fig. 7 and the video in supplementary material.

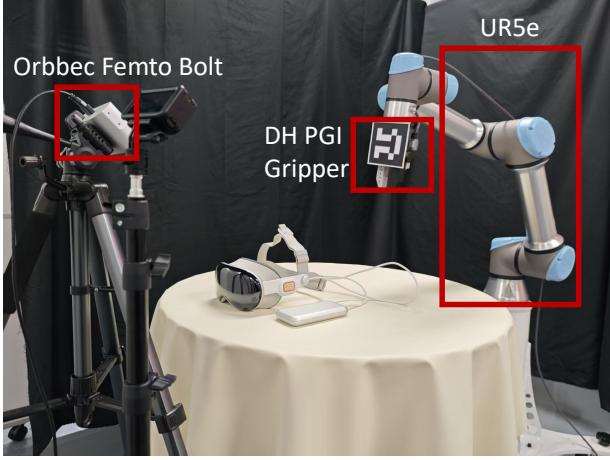


Figure 4. The experiment setting of single table-top robotic arm.

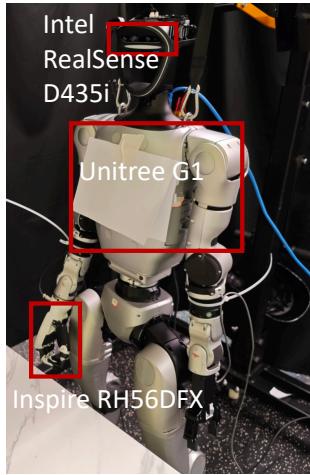


Figure 5. The experiment setting of a humanoid robot with dexterous hands.

## C. Failure Cases

The failure cases in the real-robot experiments can be summarized as two main points:

- We use the Inverse Kinematics (IK) to calculate the motions based on the target pose, which leads to rough grasping motions. The lack of dexterous grasping policy leads to collisions and overturns, as shown in the first and second row of Fig. 6. In the future, we aim to train ACT or Diffusion Policy for the grasping motion optimization.
- For the calibration of the relationship between the hand and eye of humanoid robots, our measurements may have certain errors, which could lead to the failure of tasks requiring fine operations, as shown in the third row of Fig. 7.

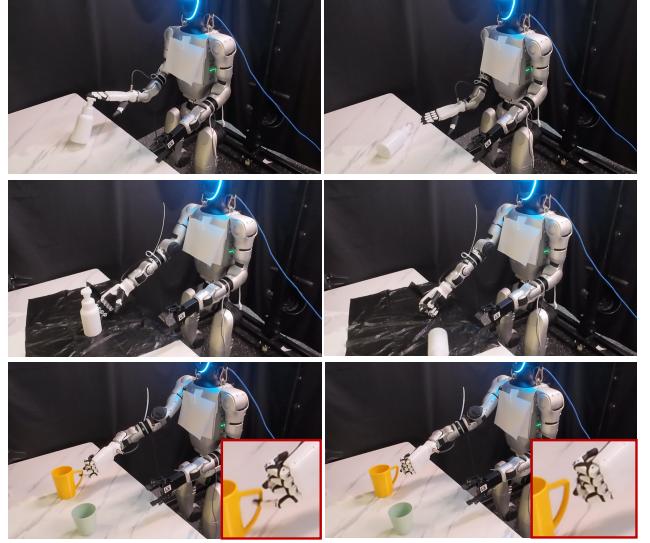


Figure 6. **Failure cases.** Row 1 & 2: Rough grasping motions lead to collisions and overturns. Row 3: The cumulative error in hand-eye calibration leads to inaccurate positions.

## D. Tracking Performance

The real-time performance of our GLOVER makes tracking moving objects while reasoning affordance possible. Time cost can be found in Table 3 in the main manuscript. We show some real-time tracking & affordance reasoning examples in Fig. 8. The qualitative results demonstrate the fast, precise, and robust tracking performance of our GLOVER. More examples can be found in the video of supplementary material.

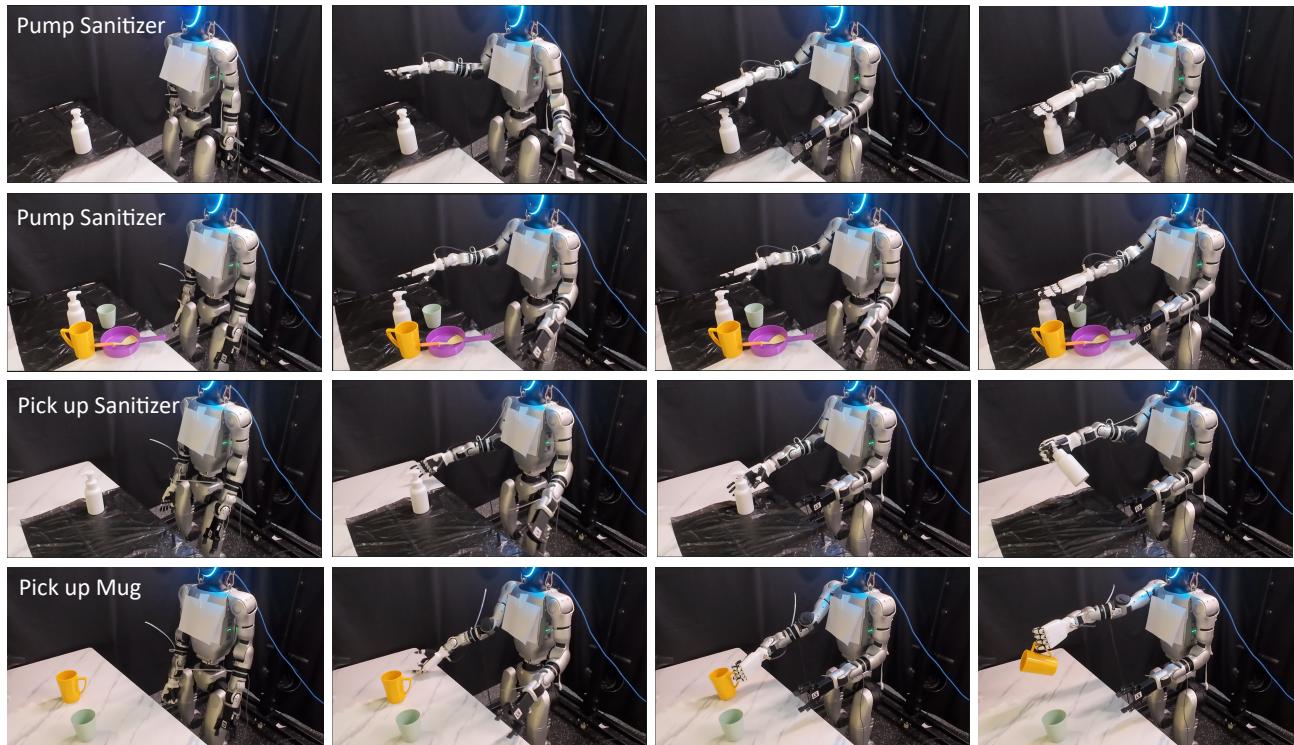
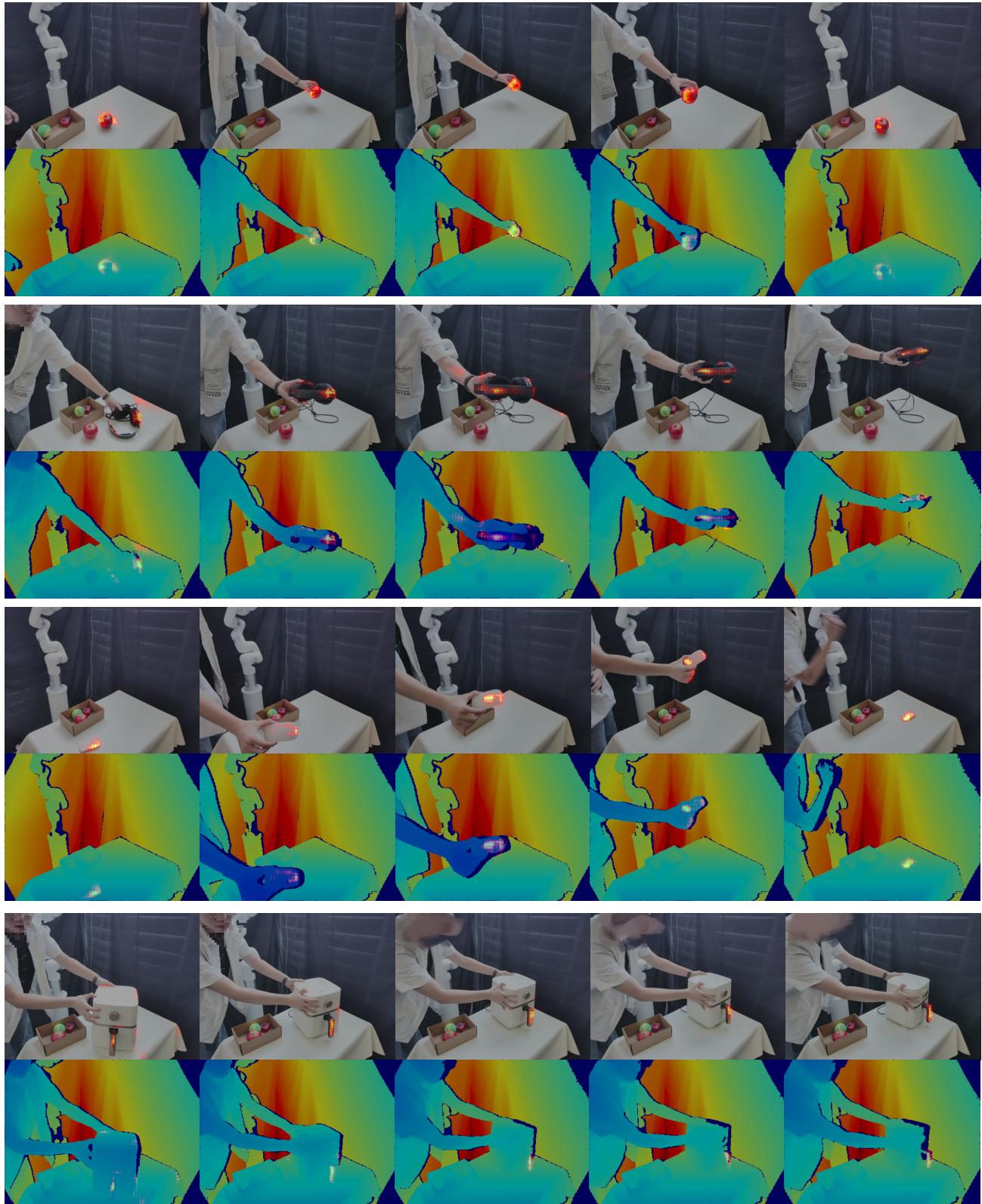


Figure 7. Demonstrations of manipulation tasks completed by the Unitree G1.



**Figure 8. The visualization of tracking performance of GLOVER.** We track the object *apple*, *earphone*, *mouse*, *oven handle* from top to bottom.