

# AffordDexGrasp: Open-set Language-guided Dexterous Grasp with Generalizable-Instructive Affordance

Yi-Lin Wei<sup>\*1</sup>, Mu Lin<sup>\*1</sup>, Yuhao Lin<sup>1</sup>, Jian-Jian Jiang<sup>1</sup>,  
Xiao-Ming Wu<sup>1</sup>, Ling-An Zeng<sup>1</sup>, Wei-Shi Zheng<sup>†1,2</sup>

<sup>1</sup> School of Computer Science and Engineering, Sun Yat-sen University, China

<sup>2</sup> Key Laboratory of Machine Intelligence and Advanced Computing, Ministry of Education, China

{weiylin5, linm67, linyh96, jiangjj35, wuxm65, zenglan3}@mail2.sysu.edu.cn wszheng@ieee.org

<https://isee-laboratory.github.io/AffordDexGrasp/>

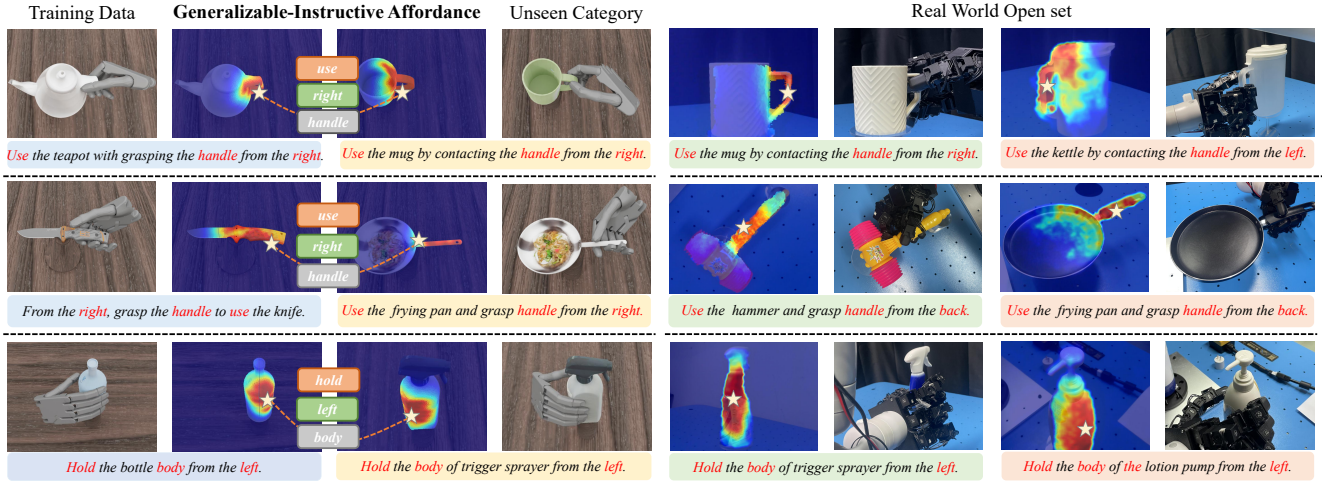


Figure 1. Open-set Language-guided Dexterous Grasp. Our framework bridges the gap between language and grasp actions through Generalizable-Instructive Affordance, which enables cross-category generalization via category-agnostic cues and graspable local structure. Remarkably, our framework demonstrates strong generalization without requiring extra real training data in real-world experiments.

## Abstract

Language-guided robot dexterous generation enables robots to grasp and manipulate objects based on human commands. However, previous data-driven methods are hard to understand intention and execute grasping with unseen categories in the open set. In this work, we explore a new task, Open-set Language-guided Dexterous Grasp, and find that the main challenge is the huge gap between high-level human language semantics and low-level robot actions. To solve this problem, we propose an Affordance Dexterous Grasp (AffordDexGrasp) framework, with the insight of bridging the gap with a new generalizable-instructive affordance representation. This affordance can generalize to unseen categories by leveraging the object’s local structure and category-agnostic semantic attributes, thereby ef-

fectively guiding dexterous grasp generation. Built upon the affordance, our framework introduces Affordance Flow Matching (AFM) for affordance generation with language as input, and Grasp Flow Matching (GFM) for generating dexterous grasp with affordance as input. To evaluate our framework, we build an open-set table-top language-guided dexterous grasp dataset. Extensive experiments in the simulation and real worlds show that our framework surpasses all previous methods in open-set generalization.

## 1. Introduction

Achieving generalizable robot dexterous grasping is an important goal in the fields of robotics and computer vision, with exciting potential applications in human-robot interaction and robot manipulation.

Recent works explore the task of language-guided dexterous grasp generation[20, 24, 25, 69], aiming to enable

<sup>\*</sup>Equal contribution.

<sup>†</sup>Corresponding author.

dexterous hands to perform actions based on language instructions, going beyond previous works that focus on stable grasping [31, 44, 54]. The typical language-guided approaches employ language as the condition of generative models to predict hand parameters, achieving impressive performance on objects within known categories [4, 59]. However, in the open real world, there are many categories that may not appear during training, and the cost of collecting data for dexterous hands is quite expensive. Therefore, the open-set generalization on unseen category samples is crucial for robot grasp. While previous works on parallel grippers explore open-set tasks [27, 42, 70], there is limited research on open-set generalization for dexterous hands, as their significantly higher degrees of freedom present complex challenges.

In this work, we explore a novel and challenging task: open-set language-guided dexterous grasping, as shown in Figure 1, where models are evaluated on objects and language instructions from both seen and unseen categories. This task poses a great challenge in ensuring that grasps are intention consistent with corresponding language instructions for unseen categories. We find that the main challenge lies in the huge gap between high-level natural language and low-level robot action spaces, which makes it difficult to generalize the ability of understanding intentions and grasping from the training domain to unseen categories.

To solve the above challenges, we propose the Affordance Dexterous Grasp (AffordDexGrasp) framework, with the insight of employing a new affordance as an intermediate representation to bridge the gap between high-level language and low-level grasp actions. Our affordance representation combines two characteristics: **generalizable** (generalizing to unseen categories based on language) and **instructive** (guiding dexterous grasp generation effectively). However, achieving these two characteristics is not trivial. For example, as shown in Figure 2, fine-grained contact information can effectively guide grasp generation or optimization [2, 9, 19, 21, 64], but it is difficult to generalize to unseen categories. In contrast, coarse information, such as object parts, can be obtained from a off-the-shelf model [32, 48, 50], but it is too coarse to guide dexterous hand actions with higher degrees of freedom.

To achieve these characteristics, we propose the Generalizable-Instructive Affordance by defining a general dexterous affordance that aligns with category-independent information, such as intention, object parts, and direction. As shown in Figure 2, the affordance represents the potential graspable regions of all grasps with the same semantics. In this way, the models do not need to learn complex dexterous contact patterns but instead focus on a general graspable area, which can be well aligned with category-agnostic semantic attributes and guide grasp generation effectively.

To generate affordance and employ it to guide grasp gen-

eration, our framework consists of two cascaded generative models. The Affordance Flow Matching generates affordance maps in a generalizable manner based on language. And the Grasp Flow Matching generates dexterous grasp poses under the effective guidance of affordance. Moreover, we introduce a pre-understanding stage and a post-optimization stage to further boost generalization. Specifically, we employ the Multimodal Large Language Model (MLLM) to pre-understand user intention to enhance generalization across diverse user commands. And we introduce an affordance-guided optimization to improve grasp quality while preserving consistency with user intentions.

For evaluating our framework, we build a open-set tabletop language-guided dexterous grasp dataset, based on the language guided dexterous grasp dataset [59, 65]. We exclude specific categories from the training set to test the model’s open-set generalization. Moreover, we provide high-quality rendered images to facilitate the usage by MLLM, and we extend the dataset to scene-level data to better simulate real-world environments. The comprehensive experiments are conducted on both the simulation open-set dataset and real-world environments. The results show that our framework can generate dexterous grasp with consistent intention and high quality in open set.

## 2. Related work

### 2.1. Dexterous Grasp

Dexterous grasping is critical in robotics, which equips robots with human-like grasping capabilities [29, 61]. Some methods [34, 51, 56, 68] focus on grasp stability and quality, while recent studies explore task-oriented [58, 72] or language-guided grasping [25, 59] with specific semantic intention. For language-guided task, the data-driven methods achieve impressive performance in unseen object within seen categories by leveraging language conditioned generative model [4, 59]. However, we find that these models are hard to generalize for samples from unseen categories. And we find that the main challenge is the huge gap between language and grasp action. To solve this problems, we propose AffordDexGrasp framework with the insight that bridge the gap by generalizable-instructive affordance.

### 2.2. Open-set Robot Grasp

Exploring the performance of robotic models in open-set scenarios is crucial due to the diverse object categories in the real world and the high cost of data collection. For tasks that only focus on stable grasping, both parallel grippers [6, 52, 60] and dexterous hands [55, 68] achieve impressive performance on open set. However, when it comes to considering task-oriented [53] or language-guided grasping [38], it becomes much more difficult. Some works in parallel grippers achieve this by using pre-trained visual ground-

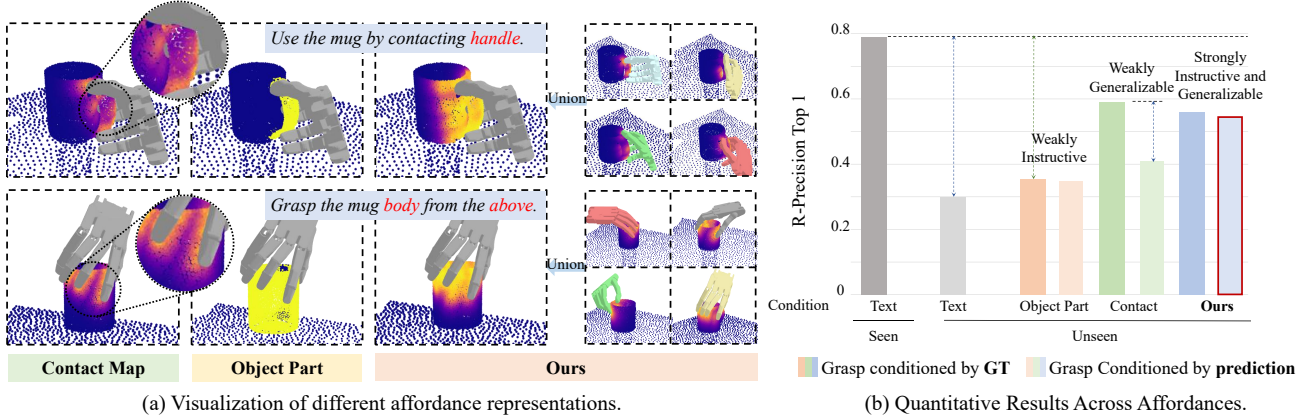


Figure 2. Different affordance representations. (a) While contact map are too elaborate to generalize and object part are too coarse to guide grasping, our affordance achieve a balance. (b) While object part has a lower upper bound and contact shows significant degradation in generalization, only our affordance effectively achieve the balance (Top-1 indicates grasp intention consistency).

ing models [26, 32, 48], implicit feature semantic fields [22, 37, 42], or multimodal large models [14, 15]. Compared to parallel grippers, dexterous hands have a significantly higher number of degrees of freedom, making the methods of parallel grippers difficult to transfer to dexterous hand [59]. In this paper, we explore the open-set language-guided dexterous grasping task and propose a novel framework to address it.

### 2.3. Grasp Affordance

Affordance is first introduced in [8], referring to environmental action possibilities. There are currently two types of grasp affordance: fine-grained contact [3, 49] and coarse-grained object segmentation [32]. The fine-grained information, such as contact areas, normals, or key points, is commonly used in hand-object interaction [2, 9, 66, 69, 73]. However, this fine-grained information is difficult to generalize well to unseen categories. As a result, inaccurate contact information would lead to unreasonable grasping. On the other hand, We find through experiments that the coarse object segmentation [32] is too coarse to guide dexterous grasp generation with relatively high degrees of freedom. To solve this problem, we propose Generalizable-Instructive Affordance, which can generalize through the local structure of objects using category-independent clues and effectively guide grasp generation.

## 3. Affordance Dexterous Grasp Framework

Given the scene point cloud  $\mathcal{O}$ , RGB images  $\mathcal{I}$  and language command  $\mathcal{C}$  as input, our goal is to generate intention aligned and high-quality dexterous grasps poses  $\mathcal{G}^{dex} = (r, t, q)$ , where  $r$  denotes the rotation,  $t$  denotes the translation, and  $q$  represents the joint angles of the dexterous hand.

In this section, we first introduce a novel generalizable-

instructive affordance representation (Section 3.1). Built on this affordance representation, we then propose the Affordance Dexterous Grasp framework (Section 3.2) including intention pre-understanding stage, Affordance and Grasp Flow Matching, and grasp post-optimization stage.

### 3.1. Generalizable-Instructive Affordance

We propose the generalizable-instructive affordance as an intermediate representation to bridge the gap between high-level language and low-level grasp actions. The affordances serves two primary objectives: 1) enable generalization by aligning object local structures with category-agnostic language semantics; and 2) provide instructive guidance for grasp generation through object affordance cues. However, there exists a trade-off in the generalization and instruction, and it is not trivial to achieve both objectives.

**Generalization vs. Instruction.** While more fine-grained contact information can better guide grasp generation, it becomes harder to generalize to unseen categories. Fine-grained affordances typically compute the distance between object points and predefined elements, such as hand surface points or hand key points [23, 33]. However, we observe that these fine-gained representations are difficult to generalize to unseen category samples as shown in Figure 2. As the dexterous hand structure is complex and contact modes are varied, the contact maps are also elaborate and varies in space, which make it difficult to generalize.

On the other hand, the information agnostic to dexterous hands, like object part segmentation, can be obtained in a generalizable manner by vision models [32, 71] or semantic feature fields [42]. However, the object part is too coarse for dexterous hand, as it may struggle to define the clear grasping area. For example, grasping the body of a mug from above and from the side involves the same part, but with different intentions. Moreover, this coarse-grained



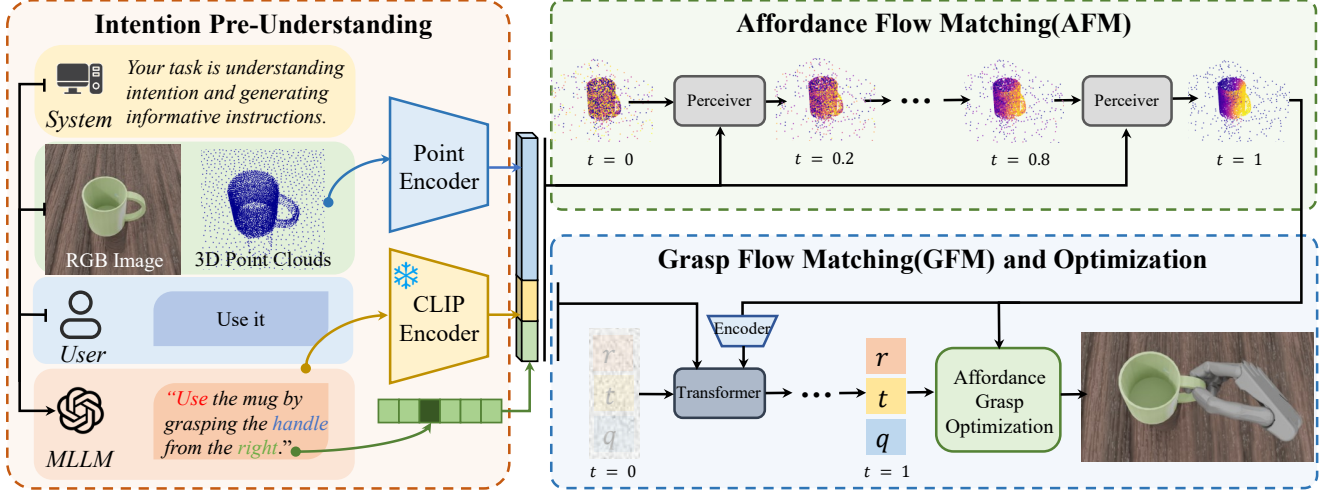


Figure 3. The pipeline of Affordance Dexterous Grasp framework. The inference pipeline includes three stages: 1) intention pre-understanding assisted by MLLM; 2) affordance flow matching for generating affordance base on MLLM output; 3) Grasp Flow Matching and Optimization for outputting grasp poses based on the affordance and MLLM outputs. In the training time, AFM and GFM are independently trained one after another. Transformer and Perceiver are attention-based interaction module for velocity vector field prediction.

learning may hinder the model from learning an accurate grasping method during training.

**Generalizable-Instructive Affordance.** To address this challenge, we propose a novel affordance representation that simultaneously enables both generalization and instruction. The key insight is to establish a correspondence between the general graspable affordance of objects and category-agnostic semantic attributes. Unlike contact maps that describe a particular hand contact area, our affordance represents all potential graspable regions that share the same semantic intention. This approach allows the affordance to generalize more effectively with category-agnostic semantic attributes from language, as it doesn’t have to learn the complex distribution of hand contact. Meanwhile, our affordance serves as priors to enhance grasp generation by providing valuable graspable area cues.

To obtain the ground truth of generalizable-instructive affordance, we first organize those grasps of one object with the same intention, contact parts and grasp direction into a semantic group. Then we calculate the distances between each scene points  $\{o_i\}^M$  and hand surface points  $\{h_i\}^N$ , to obtain the contact map  $\Omega_{contact} = \{d_i\}_1^M$  for each grasp in the group, where  $d_i = \min_{j=1 \dots N} \|o_i - h_j\|_2$ . Then we unite contact areas of  $k$  grasps in each group by calculating the minimum value of contact areas:  $d_i^{min} = \{\min_k(d_i^k)\}^M$ . In order to make the affordance map more smooth, we apply Gaussian filtering to the union contact map:  $a_i = \sum_{j \in \mathcal{N}(i)} \hat{w}_{ij} d_j^{min}$  and the weights are obtained by:

$$\hat{w}_{ij} = \frac{\exp(-\|o_i - o_j\|^2 / (2\sigma^2))}{\sum_{k \in \mathcal{N}(i)} \exp(-\|o_i - o_k\|^2 / (2\sigma^2))}, \quad (1)$$

where  $\mathcal{N}(i)$  is the neighborhood of  $o_i$ ,  $\sigma$  is set as the average nearest neighbor distance. Finally, we obtain the ground truth of generalizable-instructive affordance  $\Omega = \{a_i\}^M$ .

## 3.2. AffordDexGrasp Framework

### 3.2.1. Intention Pre-Understanding

The variability of user expression in the open world presents a challenge for models in generalizing and understanding user intentions. To address this, we employ the Multimodal Large Language Model (MLLM) to pre-understand user intentions, by inputting the user language and the rendered RGB image into GPT-4o [16] for intention understanding and key cues extraction. We empirically find that the following cues are both helpful for grasping and easily reasoned by the MLLM: the object category, user intention (e.g., use it), the contacting part of the object (e.g., the handle of a mug), and the grasp direction. The key cues can be extracted from language when the commands are clear, otherwise can be reasoned from the image.

To improve the understanding ability of MLLM, we employ chain of thought (CoT) and visual prompt, and the prompts can be found in supplementary material A.4. Additionally, the grasp direction is simplified as a discretized direction in six predefined coordinates (front, back, left, right, up, down). We first ask MLLM to obtain the direction in image coordinates, and transform it to world coordinates by the camera pose. Then the direction vector is derived from the index of the nearest coordinate axis to form discretized direction. Finally, MLLM organizes this information into a concise sentence (e.g., “use the mug from the left by contacting the handle”). This compact and information-dense



sentence structure allows subsequent models to better understand users' intentions.

### 3.2.2. Affordance Flow Matching

We introduce Affordance Flow Matching (AFM) to learn the affordance distribution efficiently, achieving intention aligned and generalizable affordance generation. AFM is built on the conditional flow matching model [1, 30, 35] aiming to learn the velocity vector field from the noised affordance map  $\Omega_0$  to the target affordance map  $\Omega_1$ , and the objective loss of AFM could be:

$$\mathcal{L}_{AFM} = \mathbb{E}_{t, \Omega_0 \sim p_0, \Omega_1 \sim p_1} \|v_\eta(\Omega_t, t|c) - (\Omega_1 - \Omega_0)\|^2, \quad (2)$$

where  $p_0$  is Gaussian distribution with zero mean and unit variance at time  $t = 0$ , and  $p_1$  is the target distribution as time  $t = 1$ .  $v_\eta(\Omega_t, t|c)$  is parameterized as a neural network with weights  $\eta$  and  $c$  is the input condition features.

Specifically, the input scene object point cloud  $\mathcal{O}$  is encoded by Pointnet++ [40], language is encoded by a pre-trained CLIP model [41], and the direction vector and time embedding are encoded by the Multilayer Perceptrons (MLP). Due to the permutation invariance of point level affordance, the noisy affordance are concatenated with object features to aware the position and semantic information. Then all features are fed into a Perceiver IO Transformer, following [17, 57]. The noisy affordance features are served as input and output array, and the features of language and direction are served as latent array. The outputs of the decoder go through an MLP to obtain the prediction of the flow  $v_\eta(\Omega_t, t|c) \in R^M$ . Finally, in the inference time, we adopt multistep  $1/\Delta t$  following the equation:

$$\hat{\Omega}_{t+\Delta t} = \Omega_t + \Delta t v_\eta(\Omega_t, t|c_1), \quad t \in [0, 1). \quad (3)$$

### 3.2.3. Grasp Flow Matching

We introduce Grasp Flow Matching (GFM) to learn the grasp distribution, conditioned on the output of MLLM and AFM. GFM is also built upon the flow matching model, and the objective is to generate a parameterized dexterous hand pose  $\mathcal{G}^{dex} = (r, t, q)$ , which represents the global rotation, translation, and joint angles. In the training time, grasp generation requires calculating explicit losses in 3D space, which involves obtaining the target dexterous poses, which go through forward kinematics to obtain hand configuration [62]. To obtain the target pose, we estimate it in one step during the training stage:

$$\mathcal{G}_1^{dex} = \mathcal{G}_t^{dex} + (1 - t) * v_\theta(\mathcal{G}_t^{dex}, t|c_2), \quad t \in [0, 1), \quad (4)$$

where the condition features  $c_2$  are the concatenation of object affordance, language and direction features. And the object affordance features are obtained by feeding affordance with point clouds into PointNet++ [40].

The losses of grasp generation include 1) the grasp pose regression L2 loss; 2) the hand chamfer loss to minimize the discrepancies between the actual hand shapes; 3) the hand fingertip key point loss to minimize the distance between fingertip position. The loss function can be formulated as:

$$\mathcal{L}_{GFM} = \lambda_{pose}\mathcal{L}_{pose} + \lambda_{chamfer}\mathcal{L}_{chamfer} + \lambda_{tip}\mathcal{L}_{tip}, \quad (5)$$

where  $\lambda_{pose}$ ,  $\lambda_{chamfer}$  and  $\lambda_{tip}$  are loss weights.

The object penetration loss, which is used to penalize the penetration of hands into objects, is excluded from grasp generation training, as it leads to significant challenges in the training process [59]. Employing the penetration loss can reduce object penetration, but negatively impacts both intention alignment and generation diversity. Therefore, to achieve semantic-aligned grasp generation, we exclude object penetration loss and introduce affordance-guided grasp optimization to reduce object penetration.

### 3.2.4. Affordance-guided Grasp Optimization

We introduce a non-parametric test-time adaptation (TTA), affordance-guided grasp optimization, to improve the grasp quality while maintaining intention consistency. Compared to learning-based TTA methods [21, 59], our non-parametric optimization does not suffer from performance degradation caused by out-of-domain open-set samples.

We design optimization objectives based on that high-quality grasps should exhibit high grasp stability and intention alignment in grasping posture and contact areas. Specifically, our optimization objectives includes: 1) affordance contact loss for constraining hand contact candidate points in contact with the affordance area:

$$\mathcal{L}_{aff-dist} = \sum_i \mathbb{I}((d(p_i^c) < \tau_1) \vee (d(p_i^r) < \tau_1)) \cdot d(p_i^r), \quad (6)$$

where  $d(p) = \min_{o \in \{o_j \in \mathcal{O} | a_j > \tau_2\}} \|p - o_j\|_2$ , and  $o_j$  and  $a_j$  represent the position and affordance values of object points, and  $p_i^c$  and  $p_i^r$  represent the hand contact candidates of the initial coarse hand and current refined hand. 2) affordance fingertip loss to keep the fingertip positions that are in contact with the affordance area unchanged.

3) object penetration loss for punish hand-object penetration, 4) self-penetration loss to punish the penetration of hand it self. 5) joint limitation loss to punish out-of-limit joint angles. The optimization objectives can be formulated as:

$$\min_{\mathcal{G}^{dex}} (\lambda_1 \mathcal{L}_{aff-dist} + \lambda_2 \mathcal{L}_{aff-finger} + \lambda_3 \mathcal{L}_{pen} + \lambda_4 \mathcal{L}_{spen} + \lambda_5 \mathcal{L}_{joint}). \quad (7)$$

## 4. Open-set Dexterous Grasp Dataset

To support the evaluation of our framework, we build an open-set dexterous grasp dataset in table-top environment based on the language-guided dexterous grasp dataset

	Intention					Quality			Diversity		
	$FID \downarrow$	$CD \downarrow$	$\frac{R - Precision \uparrow}{Top1 \quad Top2 \quad Top3}$			$Suc. \uparrow$	$Q1 \uparrow$	$Pen. \downarrow$	$\delta_t \uparrow$	$\delta_r \uparrow$	$\delta_q \uparrow$
<b>Open Set A</b>											
ContactGen[33]	0.428	9.38	0.164	0.232	0.289	11.6%	1.71e-4	1.09	2.86	18.4	39.2
Contact2Grasp[23]	0.426	10.2	0.257	0.337	0.398	16.5%	0.0172	0.668	5.73	4.84	48.3
GraspCVAE[46]	0.378	5.55	0.309	0.402	0.469	21.9%	0.0150	0.709	4.36	4.60	21.2
SceneDiffuser[13]	0.303	5.02	0.397	0.473	0.523	29.1%	0.0151	0.487	7.37	<b>9.44</b>	<b>68.9</b>
DexGYS[59]	0.297	4.63	0.317	0.401	0.463	44.2%	0.0512	0.362	6.30	6.79	54.7
Ours	<b>0.231</b>	<b>3.81</b>	<b>0.480</b>	<b>0.588</b>	<b>0.666</b>	<b>45.1%</b>	<b>0.0531</b>	<b>0.293</b>	<b>7.48</b>	6.98	61.5
<b>Open Set B</b>											
ContactGen[33]	0.492	8.95	0.196	0.257	0.299	6.04%	1.30e-4	1.11	3.26	18.7	45.5
Contact2Grasp[23]	0.369	11.9	0.248	0.318	0.367	16.2%	0.00801	0.798	5.22	4.63	40.7
GraspCVAE[46]	0.365	5.35	0.297	0.357	0.403	24.2%	0.00421	0.738	3.96	5.47	25.6
SceneDiffuser[13]	0.271	6.50	0.302	0.355	0.393	30.7%	0.00851	0.683	5.41	6.13	<b>72.9</b>
DexGYS[59]	0.358	3.40	0.294	0.358	0.403	35.2%	0.0220	0.691	5.40	6.11	59.0
Ours	<b>0.162</b>	<b>2.95</b>	<b>0.532</b>	<b>0.609</b>	<b>0.655</b>	<b>38.9%</b>	<b>0.0352</b>	<b>0.361</b>	<b>6.86</b>	<b>6.84</b>	56.3

Table 1. Results on open-set datasets compared with the SOTA methods.

[59, 65]. Our open-set dataset consists of 33 categories, 1536 objects, 1909 scenes, and 43,504 dexterous grasps for Shadow Hand [43] and Leap Hand [45].

- **Open Set Data Split.** To enable comprehensive evaluation, we construct Open sets A and B using two independent dataset splits, each with its own set of unseen categories. For each split, all categories are first labeled as either seen or unseen. Then, 80% of the objects from the seen categories are used for training, while the remaining 20% and all objects from the unseen categories are used for testing.

- **Scene Construction.** To make our dataset more practical, we generate scene data by placing objects in a table-top environment using Blenderproc [5]. To prevent collisions between the hand and the table, objects are elevated using a shelf, and collision detection is performed to filter out invalid grasps. The physics engine is activated to ensure physically plausible object placements and grasping poses, enhancing dataset quality.

- **Scene Point-Cloud and Image Rendering.** The scene was captured using five RGB-D cameras: four positioned at elevated lateral angles and one directly above the object. The partial point clouds are merged into a global point cloud. To render realistic RGB images for MLLM processing, we use the texture generation model Paint3D [67] to apply realistic textures to all objects in our dataset.

## 5. Experiment

### 5.1. Evaluation Metrics

We employ three types of metric to evaluate the ability of intention consistency, grasp quality and grasp diversity.

- **Evaluating Intention Consistency.** We employ several metrics, including FID, R-Precision and Chamfer Distance.

1) **FID** (Frechet Inception Distance) [11] measures the distribution similarity between generated grasps and ground truth. To extract grasp and instruction features, we train an object-grasp point cloud encoder and a language encoder by contrastive learning [10]. 2) **R-Precision** evaluates the semantic alignment between language instructions and generated grasps [10]. Specifically, for each generated grasp, we construct a pool of 32 samples, including its ground-truth instruction and randomly selected samples. We then compute and rank the cosine distances between the point cloud features and the language features of all samples in the pool. The average accuracy is computed at the top-1, top-2, and top-3 positions, that a retrieval is considered successful if the ground truth entry appears among the top-k candidates. 3) **Chamfer Distance (CD)** quantifies the discrepancy between the predicted hand point cloud and the closest target with same intention [59].

- **Evaluating Grasp Quality.** Following [54], we use **success rate**, denoted as  $Suc.$ , in Issac Gym [28] and **Q1** [7] to assess grasp stability. **Maximal penetration depth** (cm), denoted as  $Pen.$  is used to calculate the maximal penetration depth from the object point cloud to hand meshes.

- **Evaluating Grasp Diversity.** We compute the **Standard deviation** of translation  $\delta_t$ , rotation  $\delta_r$  and joint angle  $\delta_q$  of samples within the same scene observation.

### 5.2. Implementation Details

For training our framework, the number of epochs is set to 50 for AFM and 200 for GFM. For AFM, the loss weight of L2 loss is set to 1. For GFM, we set  $\lambda_{pose} = 10$ ,  $\lambda_{chamfer} = 1$ , and  $\lambda_{tip} = 2$ . The batch size is set to 16 for AFM and 64 for GFM. The initial learning rate is  $2.0 \times 10^{-4}$ , decaying

	<i>Open Set A</i>					<i>Open Set B</i>				
	<i>FID</i> ↓	<i>CD</i> ↓	<i>Top1</i> ↑	<i>Q1</i> ↑	<i>Pen.</i> ↓	<i>FID</i> ↓	<i>CD</i> ↓	<i>Top1</i> ↑	<i>Q1</i> ↑	<i>Pen.</i> ↓
<i>Necessity of Generalizable-Instructive Affordance</i>										
w/o affordance	0.338	5.09	0.369	0.0180	0.569	0.322	4.44	0.308	<b>0.0194</b>	<b>0.452</b>
w object part (pred)	0.313	5.06	0.320	0.0138	0.594	0.316	3.79	0.344	8.12e-3	0.721
w contact map (pred)	0.320	4.89	0.380	0.0131	0.605	0.286	3.25	0.408	9.25e-3	0.636
w our affordance (pred)	<b>0.242</b>	<b>3.79</b>	<b>0.480</b>	<b>0.0240</b>	<b>0.501</b>	<b>0.176</b>	<b>2.76</b>	<b>0.538</b>	0.0150	0.612
w object part (GT)	0.306	5.05	0.348	0.0198	0.523	0.271	3.79	0.353	6.68e-3	0.526
w contact map (GT)	0.113	2.01	0.598	0.0285	0.494	0.138	2.68	0.594	0.0143	0.643
w our affordance (GT)	0.169	3.52	0.509	0.0181	0.545	0.164	2.61	0.550	0.0120	0.584
<i>Effectiveness of Pre-Understanding Stage</i>										
w/o key cues extraction	0.258	3.99	0.463	0.0222	0.510	0.186	3.42	0.515	8.33e-3	0.702
w/o direction	0.254	4.07	0.432	0.0162	0.531	0.185	2.85	0.529	8.80e-3	0.679
w MLLM	<b>0.242</b>	<b>3.79</b>	<b>0.480</b>	<b>0.0240</b>	<b>0.501</b>	<b>0.176</b>	<b>2.76</b>	<b>0.538</b>	<b>0.0150</b>	<b>0.612</b>
<i>Effectiveness of Affordance-based Optimization</i>										
w/o our optimization	0.242	<b>3.79</b>	<b>0.480</b>	0.0240	0.501	0.176	<b>2.76</b>	<b>0.538</b>	0.0150	0.612
w penetration loss	0.336	7.07	0.318	0.0505	0.299	0.241	5.76	0.445	0.0147	0.719
w ContactNet [21]	0.349	12.3	0.279	0.0434	<b>0.124</b>	0.355	11.1	0.240	0.0206	<b>0.129</b>
w RefineNet [59]	0.249	3.96	0.436	0.0487	0.399	0.203	2.92	0.493	0.0149	0.657
w our optimization	<b>0.231</b>	3.81	0.477	<b>0.0531</b>	0.293	<b>0.162</b>	2.95	0.532	<b>0.0350</b>	0.361

Table 2. Ablation study for our framework. The results of first two experiment groups are obtained from model outputs without optimization for an intuitive and reasonable comparison. The results in each group should be compared with *light yellow line* (our default setting) .

to  $2.0 \times 10^{-5}$  using a cosine learning rate scheduler [36]. The Adam optimizer is used with a weight decay rate of  $5.0 \times 10^{-6}$ . In the inference, the time step is set to 10 for AFM and 20 for GFM. For the optimization,  $\lambda_1 = 100$ ,  $\lambda_2 = 10$ ,  $\lambda_3 = 100$ ,  $\lambda_4 = 10$ ,  $\lambda_5 = 100$ . The number of optimization iterations is set to 200. All experiments are implemented using PyTorch on a single RTX 4090 GPU.

### 5.3. Comparison with SOTA Methods

We compare our methods with the SOTA methods, as shown in Table 1. The generic grasp methods are reproduced by concatenating the point cloud features and features of same language guidance with ours. The same encoders are employed for fair comparison, and the penetration loss is excluded to avoid learning difficulties according to [59]. The results show that our framework significantly outperforms all previous methods in terms of intention consistency and grasp quality. Our framework also achieves a reasonably high level of diversity, as excessive diversity may lead to unnatural postures. Previous language-guided methods fail to generalize well to unseen categories due to the huge gap between language and grasping actions. Similarly, contact-based methods don’t perform well, as the contact maps are difficult to generalize. Additionally, our method outperforms other methods in the close set, as shown in Table 3. Overall, the results indicate that our framework achieves strong performance in generating intention-aligned, high-quality, and diverse grasps.

	<i>FID</i> ↓	<i>CD</i> ↓	<i>Top1</i> ↑	<i>Q1</i> ↑	<i>Pen.</i> ↓
GraspCVAE	0.208	2.40	0.395	0.0100	0.771
DexGYS	0.0804	1.74	0.590	0.0551	0.397
Ours	<b>0.0286</b>	<b>1.13</b>	<b>0.779</b>	<b>0.0562</b>	<b>0.193</b>

Table 3. Experiment of close set with SOTA methods.

### 5.4. Necessity of Our Affordance

The results shown in Table 2 validate the key insight of our framework: using the generalizable-instructive affordance to bridge the gap between high-level language and low-level grasp actions. The first row presents the results without affordance, which fails to generalize to unseen categories. The subsequent two rows show the results under the condition of object parts and contact maps. Prediction refers to maps generated by models, while GT refers to using the corresponding ground truth. As discussed in Section 3.1, the results demonstrate that object parts are too coarse to guide the grasp effectively, even when using the ground truth. On the other hand, the contact map is too fine-grained to generalize to unseen categories. Only our affordance achieves a balance between generalization and instruction.

Furthermore, our affordance representation and framework demonstrate good generalization in the one-shot setting, as shown in Table 4. For the one-shot experiments, we introduce several novel categories to the test set that differ significantly from the training set and add one object from each unseen category to the training set. The results show



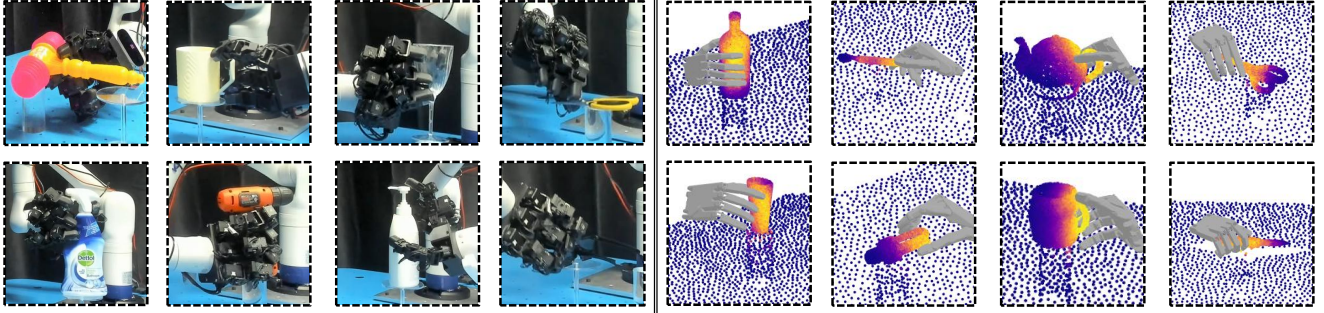


Figure 4. The visualization of generated affordance and dexterous grasp. The **left top** shows the zero-shot samples and the **left bottom** shows the one-shot samples in real world. The **right top** and **right bottom** show the zero-shot samples in simulation open set A and B.

	Seen Cate		Similar Cate		Novel Cate	
	$Top1 \uparrow$	$CD \downarrow$	$Top1 \uparrow$	$CD \downarrow$	$Top1 \uparrow$	$CD \downarrow$
<b>Zero-Shot</b>						
w/o afford	0.827	2.61	0.295	5.71	0.219	7.82
Ours	<b>0.880</b>	<b>1.06</b>	<b>0.432</b>	<b>3.65</b>	<b>0.246</b>	<b>6.42</b>
<b>One-Shot</b>						
w/o afford	0.767	3.17	0.429	3.67	0.342	6.64
Ours	<b>0.870</b>	<b>0.99</b>	<b>0.656</b>	<b>1.73</b>	<b>0.586</b>	<b>3.03</b>

Table 4. Zero-shot and one-shot generalization of our framework.

that our framework achieves a significant performance improvement, not only on samples from different categories with similar structures but also for novel categories.

### 5.5. Effectiveness of Each Component

Further ablation studies are conducted to evaluate each component of our framework. For the pre-understanding stage assisted by MLLM, *w/o key cues extraction* refers to directly inputting the user’s commands into the model, while *w/o direction* means that the direction information is not utilized. Both designs enhance performance, and we believe that MLLM would be more beneficial in real-world applications due to its powerful generalization. For affordance-based optimization, our optimization improves grasp quality and maintains intention consistency, outperforming the learning-based methods ContactNet [21] and RefineNet [59], as well as the model which is trained with penetration loss.

### 5.6. Real World Experiments

The real-world experiments are conducted to verify the simulation-to-reality ability of our framework. We employ a Leap Hand, a Kinova Gen3 6DOF arms and an original wrist RGB-D camera of Kinova arm. And we collect several common objects in daily life, including unseen objects and unseen categories. In experiment, we synthesize the scene point cloud by taking several partial depth maps





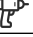



								
DexGYS	2	7	3	6	4	7	0	10
Ours	<b>10</b>	5	<b>9</b>	<b>10</b>	<b>6</b>	<b>10</b>	<b>6</b>	<b>10</b>

Table 5. Real world experiments: Successes in 10 attempts.

around the object. Then the scene point cloud, a RGB image and the user language command are fed into our framework to obtain the dexterous grasp pose. During execution, we first move the the arm to a pre-grasp position, then synchronously move the joints of the robotic arm and the dexterous hand to reach the target pose. In evaluation, a grasp is considered as success if the hand can lift it up and the action is consistent with the given command. The results shown in Table 5 shows that our framework achieve higher success rate than previous method.

## 6. Conclusions

We believe that achieving generalizable dexterous grasps in open set is crucial within the deep learning and robotics communities. In this paper, we explore a novel task of open-set language-guided dexterous grasp. This task is challenging due to the huge gap between language and grasping actions, which hinders the model’s generalization ability. We propose the AffordDexGrasp framework with the insight that using a new affordance representation, generalizable-instructive affordance, to solve this challenge. This affordance enables generalization to unseen categories by utilizing the object’s local structure and category-agnostic semantic attributes, thus facilitating effective dexterous grasp generation. Based on this affordance, our framework introduces two flow matching based models for affordance generation and affordance-guided grasp generation. Moreover, we introduce a pre-understand stage and a pose-optimization stage to further boost generalization. We conduct extensive open-set experiments in both simulation and the real world, and the results show that our framework outperforms all SOTA methods.

## 7. Acknowledgements

This work was supported partially by NSFC (92470202, U21A20471), Guangdong NSF Project (No. 2023B151504 0025), Guangdong Key Research and Development Program (No. 2024B0101040004).

## References

- [1] Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *arXiv preprint arXiv:2209.15571*, 2022. 5
- [2] Samarth Brahmabhatt, Cusuh Ham, Charles C Kemp, and James Hays. Contactdb: Analyzing and predicting grasp contact via thermal imaging. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8709–8719, 2019. 2, 3
- [3] Samarth Brahmabhatt, Chengcheng Tang, Christopher D Twigg, Charles C Kemp, and James Hays. Contactpose: A dataset of grasps with object contact and hand pose. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16*, pages 361–378. Springer, 2020. 3
- [4] Xiaoyun Chang and Yi Sun. Text2grasp: Grasp synthesis by text prompts of object grasping parts. *arXiv preprint arXiv:2404.15189*, 2024. 2
- [5] Maximilian Denninger, Dominik Winkelbauer, Martin Sundermeyer, Wout Boerdijk, Markus Wendelin Knauer, Klaus H Strobl, Matthias Humt, and Rudolph Triebel. Blenderproc2: A procedural pipeline for photorealistic rendering. *Journal of Open Source Software*, 8(82):4901, 2023. 6, 16
- [6] Hao-Shu Fang, Chenxi Wang, Hongjie Fang, Minghao Gou, Jirong Liu, Hengxu Yan, Wenhai Liu, Yichen Xie, and Cewu Lu. Anygrasp: Robust and efficient grasp perception in spatial and temporal domains. *IEEE Transactions on Robotics*, 39(5):3929–3945, 2023. 2
- [7] Carlo Ferrari, John Canny, et al. Planning optimal grasps. In *Proceedings., 1992 IEEE International Conference on Robotics and Automation, 1992.*, pages 2290–2295. IEEE, 1992. 6
- [8] James J Gibson. The theory of affordances:(1979). In *The people, place, and space reader*, pages 56–60. Routledge, 2014. 3
- [9] Patrick Grady, Chengcheng Tang, Christopher D Twigg, Minh Vo, Samarth Brahmabhatt, and Charles C Kemp. Contactopt: Optimizing contact to improve grasps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1471–1481, 2021. 2, 3
- [10] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and natural 3d human motions from text. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5152–5161, 2022. 6, 17
- [11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6, 17
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 19
- [13] Siyuan Huang, Zan Wang, Puhao Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song-Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16750–16761, 2023. 6, 17
- [14] 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. 3
- [15] Wenlong Huang, Chen Wang, Yunzhu Li, Ruohan Zhang, and Li Fei-Fei. Rekep: Spatio-temporal reasoning of relational keypoint constraints for robotic manipulation. *arXiv preprint arXiv:2409.01652*, 2024. 3, 19
- [16] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 4, 15
- [17] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Kopula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A general architecture for structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021. 5, 13
- [18] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In *International conference on machine learning*, pages 4651–4664. PMLR, 2021. 13
- [19] Juntao Jian, Xiuping Liu, Manyi Li, Ruizhen Hu, and Jian Liu. Affordpose: A large-scale dataset of hand-object interactions with affordance-driven hand pose. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14713–14724, 2023. 2
- [20] Juntao Jian, Xiuping Liu, Zixuan Chen, Manyi Li, Jian Liu, and Ruizhen Hu. G-dexgrasp: Generalizable dexterous grasping synthesis via part-aware prior retrieval and prior-assisted generation. *arXiv preprint arXiv:2503.19457*, 2025. 1
- [21] Hanwen Jiang, Shaowei Liu, Jiashun Wang, and Xiaolong Wang. Hand-object contact consistency reasoning for human grasps generation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11107–11116, 2021. 2, 5, 7, 8
- [22] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lrf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023. 3
- [23] Haoming Li, Xinzhuo Lin, Yang Zhou, Xiang Li, Yuchi Huo, Jiming Chen, and Qi Ye. Contact2grasp: 3d grasp synthesis via hand-object contact constraint. *arXiv preprint arXiv:2210.09245*, 2022. 3, 6, 17

- [24] Haosheng Li, Weixin Mao, Weipeng Deng, Chenyu Meng, Haoqiang Fan, Tiancai Wang, Ping Tan, Hongan Wang, and Xiaoming Deng. Multi-graspllm: A multimodal llm for multi-hand semantic guided grasp generation. *arXiv preprint arXiv:2412.08468*, 2024. 1
- [25] Kailin Li, Jingbo Wang, Lixin Yang, Cewu Lu, and Bo Dai. Semgrasp: Semantic grasp generation via language aligned discretization. In *European Conference on Computer Vision*, pages 109–127. Springer, 2024. 1, 2
- [26] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10965–10975, 2022. 3
- [27] 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. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10527–10534. IEEE, 2024. 2
- [28] Jacky Liang, Viktor Makoviychuk, Ankur Handa, Nuttapon Chentanez, Miles Macklin, and Dieter Fox. Gpu-accelerated robotic simulation for distributed reinforcement learning. In *Conference on Robot Learning*, pages 270–282. PMLR, 2018. 6
- [29] Yuhao Lin, Yi-Lin Wei, Haoran Liao, Mu Lin, Chengyi Xing, Hao Li, Dandan Zhang, Mark Cutkosky, and Wei-Shi Zheng. Typetele: Releasing dexterity in teleoperation by dexterous manipulation types. *arXiv preprint arXiv:2507.01857*, 2025. 2
- [30] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 5
- [31] Min Liu, Zherong Pan, Kai Xu, Kanishka Ganguly, and Dinesh Manocha. Deep differentiable grasp planner for high-dof grippers. *arXiv preprint arXiv:2002.01530*, 2020. 2
- [32] Minghua Liu, Yinhao Zhu, Hong Cai, Shizhong Han, Zhan Ling, Fatih Porikli, and Hao Su. Partslip: Low-shot part segmentation for 3d point clouds via pretrained image-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 21736–21746, 2023. 2, 3
- [33] Shaowei Liu, Yang Zhou, Jimei Yang, Saurabh Gupta, and Shenlong Wang. Contactgen: Generative contact modeling for grasp generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20609–20620, 2023. 3, 6, 17
- [34] Tengyu Liu, Zeyu Liu, Ziyuan Jiao, Yixin Zhu, and Song-Chun Zhu. Synthesizing diverse and physically stable grasps with arbitrary hand structures using differentiable force closure estimator. *IEEE Robotics and Automation Letters*, 7(1): 470–477, 2021. 2
- [35] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022. 5
- [36] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 7
- [37] Teli Ma, Zifan Wang, Jiaming Zhou, Mengmeng Wang, and Junwei Liang. Glover: Generalizable open-vocabulary affordance reasoning for task-oriented grasping. *arXiv preprint arXiv:2411.12286*, 2024. 3
- [38] 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
- [39] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 16
- [40] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 5, 13
- [41] 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. 5, 13
- [42] 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
- [43] ShadowHand. Shadowrobot. <https://www.shadowrobot.com/dexterous-hand-series/>, 2005. 6
- [44] Lin Shao, Fabio Ferreira, Mikael Jorda, Varun Nambiar, Jianlan Luo, Eugen Solowjow, Juan Aparicio Ojea, Oussama Khatib, and Jeannette Bohg. Unigrasp: Learning a unified model to grasp with multifingered robotic hands. *IEEE Robotics and Automation Letters*, 5(2):2286–2293, 2020. 2
- [45] Kenneth Shaw, Ananye Agarwal, and Deepak Pathak. Leap hand: Low-cost, efficient, and anthropomorphic hand for robot learning. *arXiv preprint arXiv:2309.06440*, 2023. 6
- [46] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems*, 28, 2015. 6, 17
- [47] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 19
- [48] Peize Sun, Shoufa Chen, Chenchen Zhu, Fanyi Xiao, Ping Luo, Saining Xie, and Zhicheng Yan. Going denser with open-vocabulary part segmentation. *arXiv preprint arXiv:2305.11173*, 2023. 2, 3
- [49] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. Grab: A dataset of whole-body human grasping of objects. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 581–600. Springer, 2020. 3



- [50] Ardian Umam, Cheng-Kun Yang, Min-Hung Chen, Jen-Hui Chuang, and Yen-Yu Lin. Partdistill: 3d shape part segmentation by vision-language model distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3470–3479, 2024. 2
- [51] Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, Yaodong Yang, Li Yi, and He Wang. Unidexgrasp++: Improving dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-specialist learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3891–3902, 2023. 2
- [52] An-Lan Wang, Nuo Chen, Kun-Yu Lin, Yuan-Ming Li, and Wei-Shi Zheng. Task-oriented 6-dof grasp pose detection in clutters. *CoRR*, abs/2502.16976, 2025. 2
- [53] An-Lan Wang, Nuo Chen, Kun-Yu Lin, Li Yuan-Ming, and Wei-Shi Zheng. Task-oriented 6-dof grasp pose detection in clutters. *arXiv preprint arXiv:2502.16976*, 2025. 2
- [54] Ruicheng Wang, Jialiang Zhang, Jiayi Chen, Yinzhen Xu, Puhao Li, Tengyu Liu, and He Wang. Dexgraspnet: A large-scale robotic dexterous grasp dataset for general objects based on simulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11359–11366. IEEE, 2023. 2, 6, 17
- [55] Wenbo Wang, Fangyun Wei, Lei Zhou, Xi Chen, Lin Luo, Xiaohan Yi, Yizhong Zhang, Yaobo Liang, Chang Xu, Yan Lu, et al. Unigrasptformer: Simplified policy distillation for scalable dexterous robotic grasping. *arXiv preprint arXiv:2412.02699*, 2024. 2
- [56] Yan-Kang Wang, Chengyi Xing, Yi-Lin Wei, Xiao-Ming Wu, and Wei-Shi Zheng. Single-view scene point cloud human grasp generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 831–841, 2024. 2
- [57] Zan Wang, Yixin Chen, Baoxiong Jia, Puhao Li, Jinlu Zhang, Jingze Zhang, Tengyu Liu, Yixin Zhu, Wei Liang, and Siyuan Huang. Move as you say interact as you can: Language-guided human motion generation with scene affordance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 433–444, 2024. 5
- [58] Wei Wei, Peng Wang, Sizhe Wang, Yongkang Luo, Wanyi Li, Daheng Li, Yayu Huang, and Haonan Duan. Learning human-like functional grasping for multi-finger hands from few demonstrations. *IEEE Transactions on Robotics*, 2024. 2
- [59] Yi-Lin Wei, Jian-Jian Jiang, Chengyi Xing, Xian-Tuo Tan, Xiao-Ming Wu, Hao Li, Mark Cutkosky, and Wei-Shi Zheng. Grasp as you say: Language-guided dexterous grasp generation. *Advances in Neural Information Processing Systems*, 37:46881–46907, 2025. 2, 3, 5, 6, 7, 8, 17
- [60] Xiao-Ming Wu, Jia-Feng Cai, Jian-Jian Jiang, Dian Zheng, Yi-Lin Wei, and Wei-Shi Zheng. An economic framework for 6-dof grasp detection. In *European Conference on Computer Vision*, pages 357–375. Springer, 2024. 2
- [61] Chengyi Xing, Hao Li, Yi-Lin Wei, Tian-Ao Ren, Tianyu Tu, Yuhao Lin, Elizabeth Schumann, Wei-Shi Zheng, and Mark R Cutkosky. Taccap: A wearable fbg-based tactile sensor for seamless human-to-robot skill transfer. *arXiv preprint arXiv:2503.01789*, 2025. 2
- [62] Guo-Hao Xu, Yi-Lin Wei, Dian Zheng, Xiao-Ming Wu, and Wei-Shi Zheng. Dexterous grasp transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 5
- [63] Yinzhen Xu, Weikang Wan, Jialiang Zhang, Haoran Liu, Zikang Shan, Hao Shen, Ruicheng Wang, Haoran Geng, Yijia Weng, Jiayi Chen, et al. Unidexgrasp: Universal robotic dexterous grasping via learning diverse proposal generation and goal-conditioned policy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4737–4746, 2023. 14
- [64] Lixin Yang, Xinyu Zhan, Kailin Li, Wenqiang Xu, Jiefeng Li, and Cewu Lu. Cpf: Learning a contact potential field to model the hand-object interaction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11097–11106, 2021. 2
- [65] Lixin Yang, Kailin Li, Xinyu Zhan, Fei Wu, Anran Xu, Liu Liu, and Cewu Lu. Oakink: A large-scale knowledge repository for understanding hand-object interaction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20953–20962, 2022. 2, 6
- [66] Ling-An Zeng, Guohong Huang, Yi-Lin Wei, Shengbo Gu, Yu-Ming Tang, Jingke Meng, and Wei-Shi Zheng. Chainhoi: Joint-based kinematic chain modeling for human-object interaction generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 12358–12369, 2025. 3
- [67] Xianfang Zeng, Xin Chen, Zhongqi Qi, Wen Liu, Zibo Zhao, Zhibin Wang, Bin Fu, Yong Liu, and Gang Yu. Paint3d: Paint anything 3d with lighting-less texture diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4252–4262, 2024. 6, 15
- [68] Jialiang Zhang, Haoran Liu, Danshi Li, XinQiang Yu, Haoran Geng, Yufei Ding, Jiayi Chen, and He Wang. Dexgraspnet 2.0: Learning generative dexterous grasping in large-scale synthetic cluttered scenes. In *8th Annual Conference on Robot Learning*, 2024. 2
- [69] Zhongqun Zhang, Hengfei Wang, Ziwei Yu, Yihua Cheng, Angela Yao, and Hyung Jin Chang. Nl2contact: Natural language guided 3d hand-object contact modeling with diffusion model. In *European Conference on Computer Vision*, pages 284–300. Springer, 2024. 1, 3
- [70] Jiaming Zhou, Ke Ye, Jiayi Liu, Teli Ma, Zifan Wang, Ronghe Qiu, Kun-Yu Lin, Zhilin Zhao, and Junwei Liang. Exploring the limits of vision-language-action manipulations in cross-task generalization. *CoRR*, abs/2505.15660, 2025. 2
- [71] Yuchen Zhou, Jiayuan Gu, Xuanlin Li, Minghua Liu, Yunhao Fang, and Hao Su. Partslip++: Enhancing low-shot 3d part segmentation via multi-view instance segmentation and maximum likelihood estimation. *arXiv preprint arXiv:2312.03015*, 2023. 3
- [72] Tianqiang Zhu, Rina Wu, Xiangbo Lin, and Yi Sun. Toward human-like grasp: Dexterous grasping via semantic representation of object-hand. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15741–15751, 2021. 2

- [73] Binghui Zuo, Zimeng Zhao, Wenqian Sun, Xiaohan Yuan, Zhipeng Yu, and Yangang Wang. Graspdiff: Grasping generation for hand-object interaction with multimodal guided diffusion. *IEEE Transactions on Visualization and Computer Graphics*, 2024. [3](#)

# Supplemental Material

## A. Framework Details

### A.1. Affordance Flow Matching

The training pipeline of Affordance Flow Matching (AFM) is shown in Figure 5 (a). A noise  $\Omega_0$  is first sampled from a Gaussian distribution with zero mean and one unit. A timestep  $t$  is sampled from 0 to unit variance. Then, the training input  $\Omega_t$  is obtained by  $\Omega_t = (1 - t) * \Omega_0 + t * \Omega_1$ , where  $\Omega_1$  is the ground truth of the affordance. The scene point clouds  $\mathcal{O}$  are fed into PointNet++ [40] with several upsampling and downsampling layers to obtain features  $\mathcal{F}_{obj} \in \mathbb{R}^{N_{scene} \times C_{scene}}$ . The object features are then concatenated with  $\Omega_t$  and point cloud positions  $\mathcal{O}$ , and processed by an MLP layer to obtain affordance features  $\mathcal{F}_{aff} \in \mathbb{R}^{N_{scene} \times C_{aff}}$ . The direction vector, a one-hot vector of 6 dimensions, is fed into an MLP to obtain the direction features  $\mathcal{F}_{dir} \in \mathbb{R}^{1 \times C_{dir}}$ . The language text is encoded by the CLIP model [41] and outputs language features  $\mathcal{F}_{lan} \in \mathbb{R}^{1 \times C_{lan}}$ . These features are then fed into Perceiver IO [17, 18] for effective feature interaction. Specifically, the direction and language features serve as the query, while the affordance features serve as the key and value for cross-attention. Then, the output of the first attention mechanism performs self-attention. Then, the scene features serve as the query and the self-attention result serve as key and value in the final cross-attention to obtain the final features. The features are then fed into an MLP to predict the flow  $v_\eta(\Omega_t, t|c_1)$ . The loss function is the L2 loss between predicted flow with the target:

$$\mathcal{L}_{AFM} = \|v_\eta(\Omega_t, t|c) - (\Omega_1 - \Omega_0)\|^2. \quad (8)$$

### A.2. Grasp Flow Matching

The training pipeline of Grasp Flow Matching (GFM) is shown in Figure 5. The object and affordance are concatenated first then fed into PointNet++ with several downsampling layer to obtain the affordance features. The direction and language are encoded same with AFM. And all features are concatenated in the token dimension to obtain the condition features. The noised grasp poses are obtained similar with AFM, which are encoded by MLP. Then the pose features serve as query, and the condition features serve as key and value in the transformer decoder. The output features are fed into an MLP to obtain the prediction pose flow. Then the grasp pose can be obtained by one step, during the training time:

$$\hat{\mathcal{G}}_1^{dex} = \mathcal{G}_t^{dex} + (1 - t) * v_\theta(\mathcal{G}_t^{dex}, t|c_2), \quad t \in [0, 1), \quad (9)$$

The loss functions includes grasp pose regression loss, hand chamfer loss and hand fingertip key point loss, which are

detailed in next section.

### A.3. Loss Function

This section provides a detailed exposition of the loss functions utilized during the grasp generation and optimization.

**- Parameter Regression Loss.** We utilize the mean squared error (MSE) to quantify the deviation between the generated dexterous hand pose  $\mathcal{G}_1^{dex}$  and the ground truth  $\mathcal{G}_1^{dex}$ .

$$\mathcal{L}_{pose} = \|\mathcal{G}_1^{dex} - \hat{\mathcal{G}}_1^{dex}\|_2^2. \quad (10)$$

**- Hand Chamfer Loss.** The predicted hand point clouds  $\mathcal{H}(\hat{\mathcal{G}}_1^{dex})$  and the ground truth  $\mathcal{H}(\mathcal{G}_1^{dex})$  are derived by sampling from the hand mesh. We then compute the chamfer distance to assess the discrepancies between the predicted and ground-truth hand shapes.

$$\begin{aligned} \mathcal{L}_{chamfer} = & \sum_{x \in \mathcal{H}(\mathcal{G}^{dex})} \min_{y \in \mathcal{H}(\hat{\mathcal{G}}_1^{dex})} \|x - y\|_2^2 \\ & + \sum_{x \in \mathcal{H}(\hat{\mathcal{G}}^{dex})} \min_{y \in \mathcal{H}(\mathcal{G}_1^{dex})} \|x - y\|_2^2. \end{aligned} \quad (11)$$

**- Fingertip Keypoint Loss.** The fingertip keypoint loss  $\mathcal{L}_{tip}$  measures the distance between the fingertip of generated grasp  $q^g$  and the ground truth  $q^{gt}$ .

$$\mathcal{L}_{tip} = \|q^g - q^{gt}\|_2^2. \quad (12)$$

**- Affordance contact loss.** The affordance contact loss is designed for constraining hand contact candidate points in contact with the affordance area:

$$\mathcal{L}_{\text{aff-dist}} = \sum_i \mathbb{I}((d(p_i^c) < \tau_1) \vee (d(p_i^r) < \tau_1)) \cdot d(p_i^r), \quad (13)$$

where  $d(p) = \min_{o \in \{o_j \in \mathcal{O} | a_j > \tau_2\}} \|p - o_j\|_2$ , and  $o_j$  and  $a_j$  represent the position and affordance values of object points, and  $p_i^c$  and  $p_i^r$  represent the hand contact candidates of the initial coarse hand and current refined hand.

**- Affordance fingertip loss.** The affordance fingertip loss is designed to keep the fingertip positions that are in contact with the affordance area unchanged.

$$\mathcal{L}_{\text{aff-finger}} = \sum_i \mathbb{I}(d(q_i^c) < \tau_3)) \cdot \|q_i^r - q_i^c\|_2 \quad (14)$$

where  $q_i^c$  and  $q_i^r$  represent the fingertip position of the initial coarse hand and current refined hand.

**- Object Penetration Loss.** The object penetration loss  $\mathcal{L}_{pen}$  penalizes the depth of hand-object penetration, where



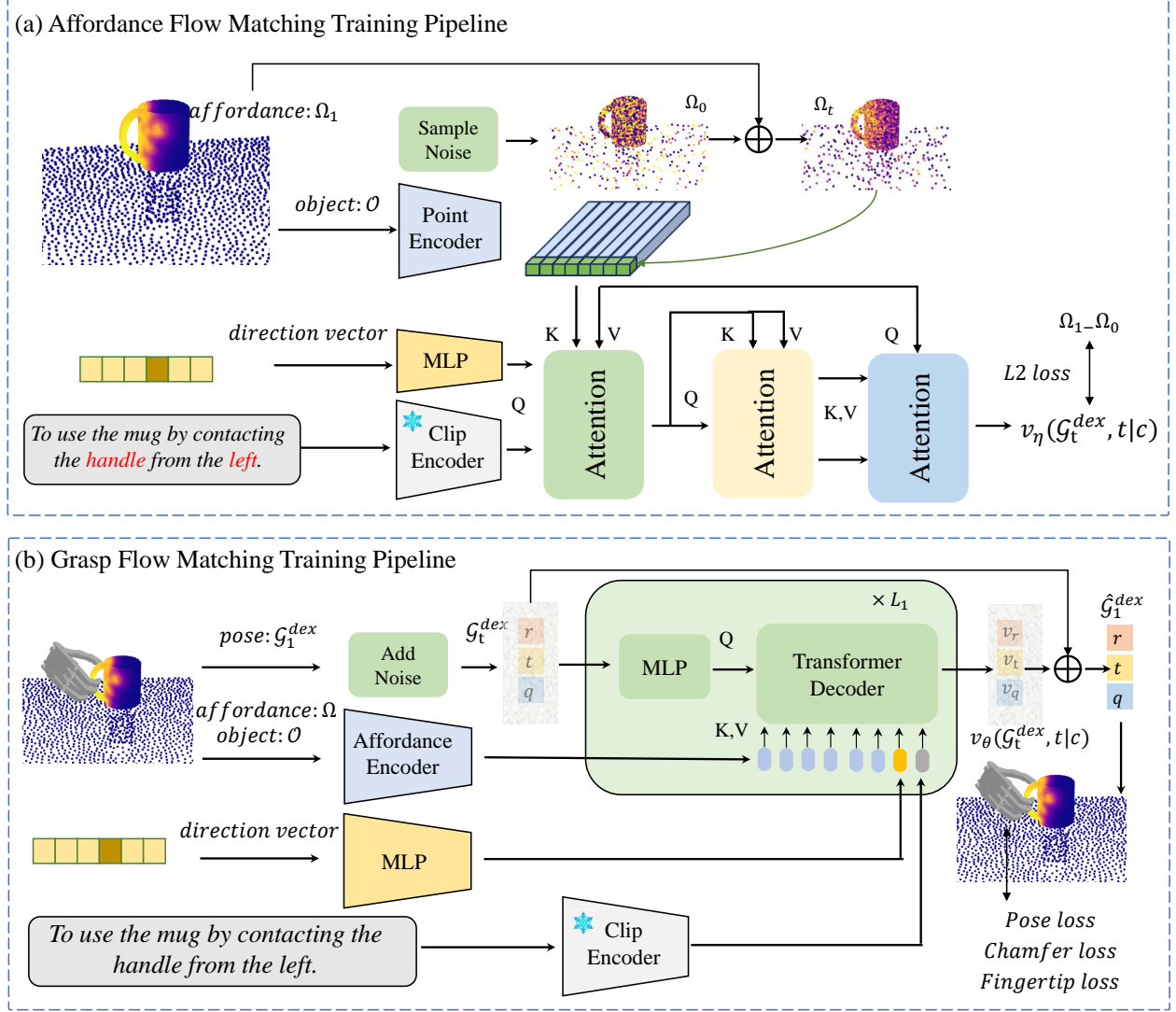


Figure 5. The training pipeline of Affordance Flow Matching and Grasp Flow Matching.

$d_i^{sdf}$  denotes the signed distance from the object point to the hand mesh.

$$\mathcal{L}_{pen} = \sum_i \mathbb{I}(d_i^{sdf} > 0) \cdot d_i^{sdf}. \quad (15)$$

**- Self-Penetration Loss.** The self-penetration loss  $\mathcal{L}_{spen}$  punishes the penetration among the different parts of the hand, where  $p^{dex,sp}$  denotes predefined anchor spheres on the hand [63].

$$\mathcal{L}_{spen} = \sum_{i,j} \mathbb{I}(i \neq j) \cdot \max(\delta - d(p_i^{dex,sp}, p_j^{dex,sp}), 0). \quad (16)$$

**- Joint Limitation Loss.** Given the physical structure limitations of the robotic hand, each joint has designated upper

and lower limits. The joint angle loss penalizes deviations from these limits.

$$\mathcal{L}_{joint} = \sum_i (\max(q_i - q_i^{max}, 0) + \max(q_i^{min} - q_i, 0)). \quad (17)$$

#### A.4. Pre-understanding by MLLM

We employ Chain-of-Thought (CoT) and visual prompt to enhance the reasoning ability of MLLM as shown in Figure 6. Specifically, the MLLM is asked to reason the following task step by step: recognize the object category, understand the intention, reason the contact part based on intention and category, reason the direction based on the above. Furthermore, we add arrows of different colors in the four direc-







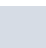

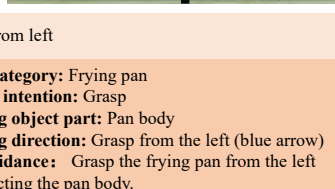


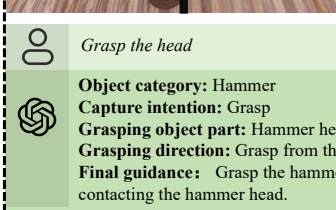
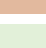
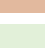
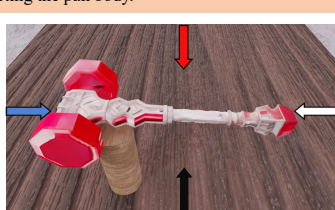
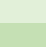
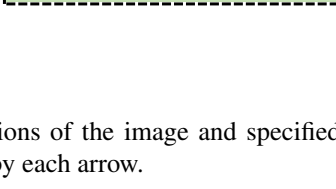
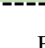
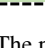
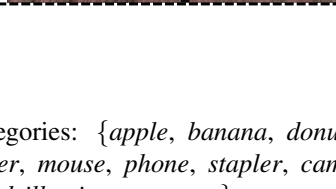
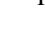
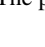
<b>System Prompt</b> : Your task is to guide the robot to execute user commands. Firstly, accept human commands, then analyze the images and output information in a specified format to guide the robot to grasp and manipulate objects.			
Please analyze step by step in the following order: 1. What category does the object in the picture belong to 2. What is the intention of grabbing (use/hold/lift/hand over) 3. What part of this object should be grasped (e.g. handle or body) 4. In which direction should this part be grasped (clockwise from red to blue indicates grasping from above, grasping from right, grasping from front, grasping from left). 5. Final guidance: <intention> the <category> from <direction> by contacting <object part>.			
And output: <b>Object category:</b> {} <b>Capture intention:</b> {} <b>Grasping object part:</b> {} <b>Grasping direction:</b> {} <b>Final guidance:</b> {}			
	 Use it	 Grasp from above	
	 <b>Object category:</b> Mug <b>Capture intention:</b> Use <b>Grasping object part:</b> Handle <b>Grasping direction:</b> Grasp from the left (blue arrow) <b>Final guidance:</b> Use the mug from the left by contacting the handle.	 <b>Object category:</b> Mug <b>Capture intention:</b> Grasp <b>Grasping object part:</b> Mug body <b>Grasping direction:</b> Grasp from above (red arrow) <b>Final guidance:</b> Grasp the mug from above by contacting the body.	
	 Grasp the right ear	 Lift it	
	 <b>Object category:</b> Headphones <b>Capture intention:</b> Grasp <b>Grasping object part:</b> Right ear cup <b>Grasping direction:</b> Grasp from the right (white arrow) <b>Final guidance:</b> Grasp the headphones from the right by contacting the right ear cup.	 <b>Object category:</b> Headphones <b>Capture intention:</b> Lift <b>Grasping object part:</b> Headband <b>Grasping direction:</b> Grasp from above (red arrow) <b>Final guidance:</b> Lift the headphones from above by contacting the headband.	
	 Use it	 Grasp from left	
	 <b>Object category:</b> Frying pan <b>Capture intention:</b> Use <b>Grasping object part:</b> Handle <b>Grasping direction:</b> Grasp from the right (white arrow) <b>Final guidance:</b> Use the frying pan from the right by contacting the handle.	 <b>Object category:</b> Frying pan <b>Capture intention:</b> Grasp <b>Grasping object part:</b> Pan body <b>Grasping direction:</b> Grasp from the left (blue arrow) <b>Final guidance:</b> Grasp the frying pan from the left by contacting the pan body.	
	 Grasp the head	 Use it	
	 <b>Object category:</b> Hammer <b>Capture intention:</b> Grasp <b>Grasping object part:</b> Hammer head <b>Grasping direction:</b> Grasp from the left (blue arrow) <b>Final guidance:</b> Grasp the hammer from the left by contacting the hammer head.	 <b>Object category:</b> Hammer <b>Capture intention:</b> Use <b>Grasping object part:</b> Handle <b>Grasping direction:</b> Grasp from the right (white arrow) <b>Final guidance:</b> Use the hammer from the right by contacting the handle.	

Figure 6. The prompt used in per-understanding stage.

tions of the image and specified the direction represented by each arrow.

## B. Zero-shot Datasets Details

### B.1. Dataset Splitting

To support our Cross-Category Open-set settings, we divided the dataset into Open Set A and Open Set B by excluding specific samples from the training set. Specifically, the test set of Open Set A consists of 9,688 samples of 808 objects from 11 unseen categories: {*bowl, cylinder, wineglass, headphones, flashlight, pen, wrench, toothbrush, hammer, teapot, scissors*}. The test set of Open Set B consists of 29,744 samples of 568 unseen objects from 10 unseen categories: {*bottle, cup, squeezable, eyeglasses, light-bulb, fryingpan, knife, screwdriver, mug, pincer*}. For one-shot experiments, we add an additional 3,688 samples of

152 objects from 12 categories: {*apple, banana, donut, binoculars, gamecontroller, mouse, phone, stapler, cameras, lotion\_pump, power drill, trigger sprayer*}.

### B.2. Texture generation

As part of the texture generation, we utilized GPT-4o[16] to generate possible colors and materials for specific object categories. The prompt used for this process was as follows:

**Prompt:** “Please generate a list of possible colors and materials for the following categories and provide the output in JSON format. The objects are: {binoculars, mug, ...}.”

For objects lacking texture, we employ Paint3D[67] to generate their visual appearance. The positive prompt is constructed using the object’s category name along with a sampled color and material from the predefined possible lists, while the default negative prompt is utilized to guide

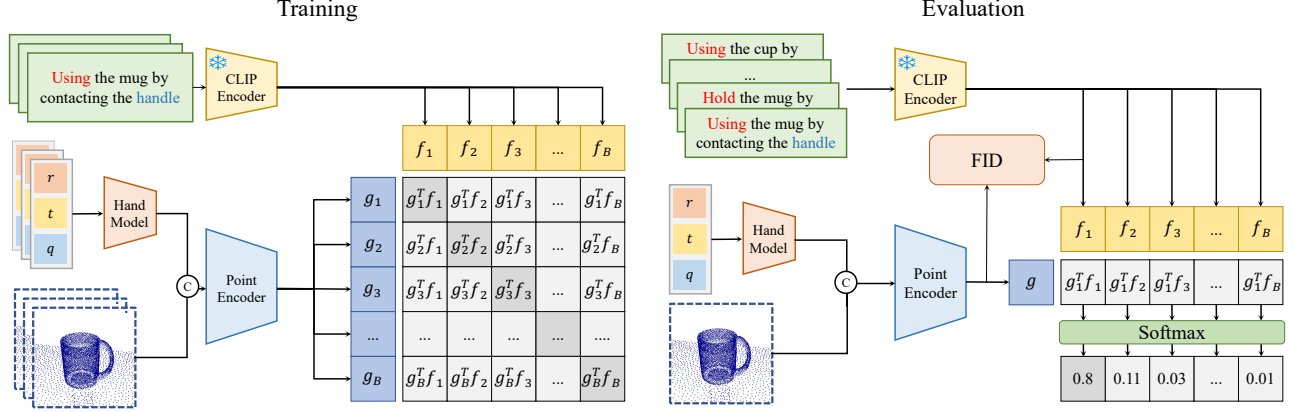


Figure 7. The training pipeline of evaluation encoder (left). The evaluation of FID and R-Precision.

the generation process.

**Prompt:** “<category>, <sampld color>, <sampld material>, high quality, clear”

**Negative prompt:** “strong light, Bright light, intense light, dazzling light, brilliant light, radiant light, Shade, darkness, silhouette, dimness, obscurity, shadow”

### B.3. Scene Reconstruction

The details of the construction of our scene dataset are illustrated in Figure 8. We utilize BlenderProc [5] to create a tabletop scene and project the grasping poses from the full-model dataset into the scene. After performing a gravity check on the objects and a collision test between the objects and the robotic hand, we record the grasp parameters in the scene coordinate system. At the same time, we capture RGB and depth images from the scene cameras. Finally, we construct the scene point cloud by concatenating and downsampling the partial point clouds obtained from the depth images.

In our tabletop scene, objects are lifted using three types of shelves. The first is a cylindrical base, which is used for the majority of categories. The second is a smaller shelf, consisting of two circular planes connected by a rod, designed for smaller categories such as pens and toothbrushes. The third is a headphone-specific shelf, used exclusively for headphones.

The poses of five RGB-D cameras are determined based on the center of the object’s bounding box. One camera is positioned directly above the center of the bounding box. The other four cameras are placed at the four sides of the bounding box, with a side offset equal to half the length of the bounding box’s longest edge and a height offset equal to one-fourth of the bounding box’s longest edge.

## C. Experiments Details

### C.1. Evaluation Details

**- Training pipeline of evaluation encoder.** The training pipeline of our Intention Evaluation Model is illustrated in Figure 7. The model takes as input the hand parameters  $\mathcal{G}^{dex}$ , the object point cloud  $\{o_i\}_{i=1}^M$ , and the language guidance  $\mathcal{L}$ . Initially, the hand model generates the hand point cloud  $\{h_i\}_{i=1}^N$ , which is then concatenated with the object point cloud  $\{o_i\}_{i=1}^M$  to form the combined hand-object point cloud  $\{p_i\}_{i=1}^P$ .

A PointNet++ encoder and a CLIP encoder are subsequently employed to map  $\{p_i\}_{i=1}^P$  and  $\mathcal{L}$  into the feature space, producing the grasping feature  $\mathbf{g}$  and the guidance feature  $\mathbf{f}$ , respectively. The Information Noise-Contrastive Estimation (InfoNCE) Loss[39] is then applied to maximize the cosine similarity between paired features while minimizing it between unpaired features. The InfoNCE Loss is defined as:

$$\mathcal{L}_{\text{InfoNCE}} = -\log \frac{\sum_{j=1}^B \mathbb{I}[\text{match}(i, j)] \exp(\text{sim}(\mathbf{g}_i, \mathbf{f}_j)/\tau)}{\sum_{j=1}^B \exp(\text{sim}(\mathbf{g}_i, \mathbf{f}_j)/\tau)}, \quad (18)$$

where  $\text{match}(i, j)$  indicates that the object category and action of the  $i$ -th grasp match those of the  $j$ -th grasp (including the case where  $i = j$ ). Here,  $\text{sim}(\mathbf{g}, \mathbf{f})$  denotes the cosine similarity between features  $\mathbf{g}$  and  $\mathbf{f}$ ,  $\tau$  is a temperature parameter, and  $B$  is the number of samples in the batch. The loss encourages high similarity for matched pairs  $(\mathbf{g}_i, \mathbf{f}_j)$  and low similarity for unmatched pairs  $(\mathbf{g}_i, \mathbf{f}_j)$  where  $i \neq j$  and  $\text{match}(i, j)$  does not hold.

To ensure the model’s ability to distinguish between different actions within the same category, we design a custom sampler for the training data loader. The sampler randomly selects 8 different categories and randomly chooses



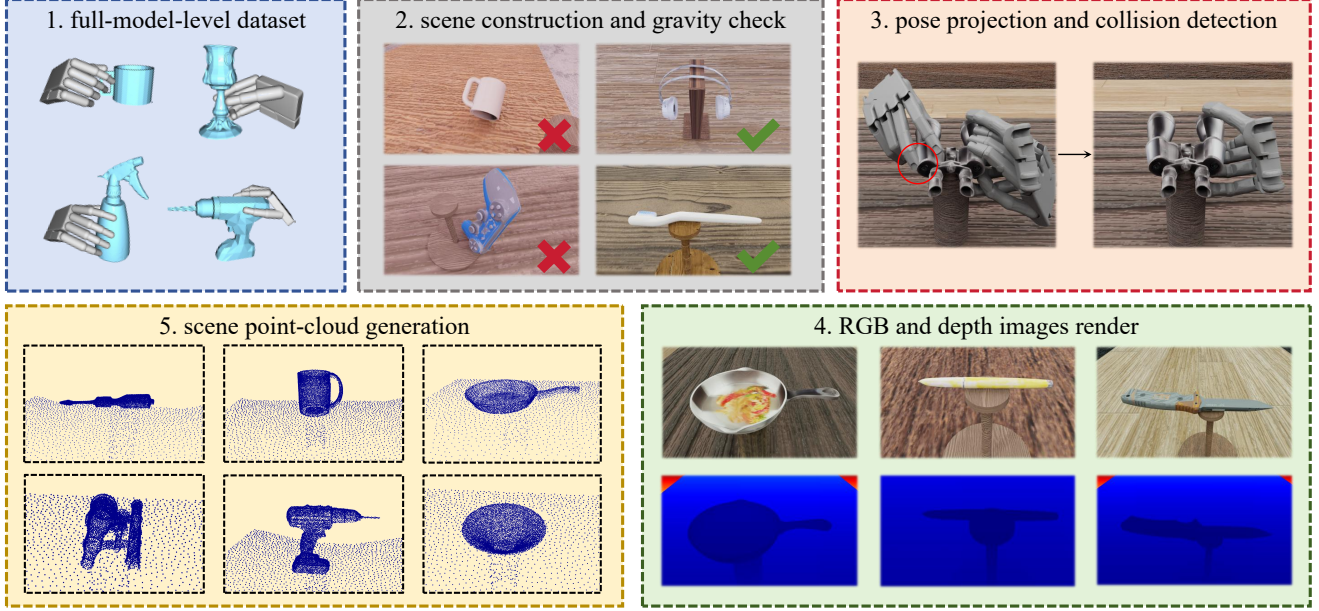


Figure 8. The construction of our scene dataset.

$\frac{B}{8}$  grasps for each category. This process is repeated until the number of remaining categories in the dataset is fewer than 8.

- **R-precision** When computing R-Precision during evaluation, we first randomly select 3 grasp guidance samples from the dataset that belong to the same category but involve different actions, along with 28 guidance samples from other categories. Next, we generate grasp features and a batch of guidance features using the same method as in the training stage. We then compute the cosine similarity between the grasp feature and each guidance feature. Finally, we determine the rank of the matched guidance sample based on these similarity scores.

- **Fréchet Inception Distance (FID)**. This metric is widely used in generative tasks [10, 11] to assess the similarity between the generated and ground truth distributions. For our evaluation, we compute the FID using features extracted from the object and hand point clouds.

- **Chamfer Distance (CD)**. Chamfer distance, denoted as  $CD$ , measures the discrepancy between the predicted hand point clouds and the nearest ground truth targets, providing a consistency measure from the perspective of hand positioning. The closest targets are the ground truth grasps that share the same intended grasp and are most similar to the generated grasp. This methodology follows the approach in [59].

- **Grasp Success Rate**. The success rate of grasps is evaluated within the Isaac Gym simulation environment. To mimic the forces exerted by a dexterous hand in real-world scenarios, each finger is contracted towards the object. A

grasp is considered successful if it can resist at least one of the six gravitational directions, indicating it can sustain a stable hold on the object.

- **Mean  $Q_1$  Metric**. The  $Q_1$  metric quantifies the smallest wrench capable of destabilizing a grasp. We follow the guidelines in [54], where the contact threshold is set at 1 cm and the penetration threshold is set at 5 mm. Any grasp with a maximal penetration depth exceeding 5 mm is considered invalid, and its  $Q_1$  value is set to 0 before averaging the results.

- **Maximal Penetration Depth**. This metric measures the greatest penetration from the object’s point cloud to the hand mesh during a grasp.

- **Diversity**. To assess the diversity of generated grasps, we calculate the standard deviations of translation ( $\delta_t$ ), rotation ( $\delta_r$ ), and joint angles ( $\delta_q$ ). We generate eight samples for each intended grasp under identical input conditions, and each sample is then sent for refinement through the quality component. For calculation purposes,  $\delta_r$  and  $\delta_q$  are converted into Euler angles (in degrees), and  $\delta_t$  is measured in centimeters.

## C.2. Reproduction of SOTA method

We reproduce SOTA methods on our Open Set dataset using the same encoder structure to ensure fair comparison. Specifically, we reimplement ContactGen [33], Contact2Grasp [23], GraspCAVE based on [46], SceneDiffuser [13], and DexGYS [59]. To introduce language information, we use an identical CLIP language encoder. For GraspCAVE, ContactGen and Contact2Grasp, we concate-

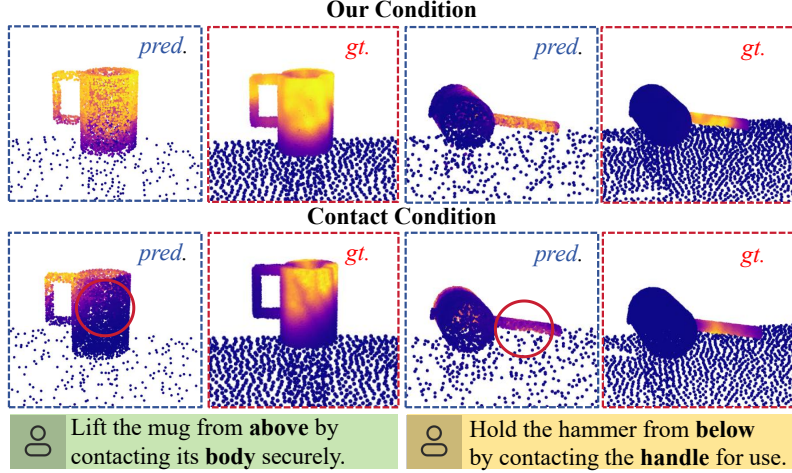


Figure 9. The visualization of our affordance and contact map for prediction and the ground truth.

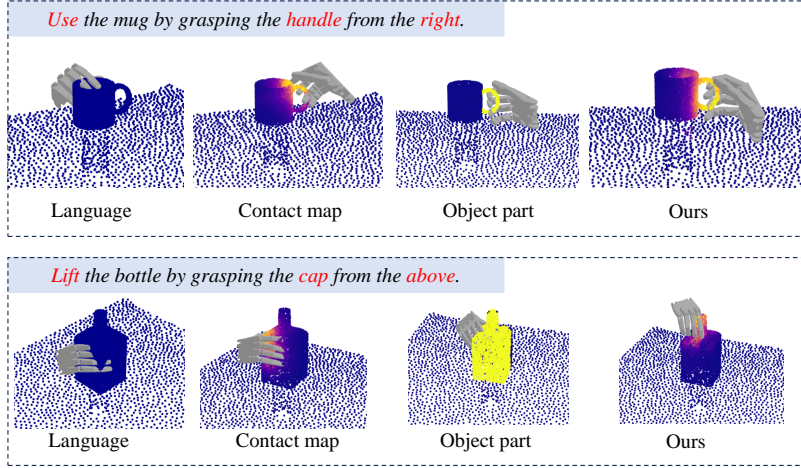


Figure 10. The visualization of generated grasps under different condition with *only language*, *contact map*, *object part* and *ours*

nate the language feature, object feature, and latent feature to send to the decoder. For SceneDiffuser, we concatenate the language and object features as the model condition. The input language guidance of all methods are identical as the MLLM’s output for fair comparison.

## D. More Experiments

### D.1. Effectiveness of Flow matching model

The results in Table 6 and 7 shows the performance of Our flow matching based model with the diffusion based model. It is noticed that our model surpass the diffusion based model with higher performance and faster inference speed. In addition, the introduced of the affordance flow matching only increased the time by 0.075 second compared to the baseline.

### D.2. Experiments of Affordance

**Qualitative Experiments.** As shown in Figure 9 and 10, our generalizable-instructive affordance is more generalizable and instructive than other representations on open set. The language condition, without low level affordance condition struggles to generalize to unseen categories. The contact map is too fine-grained to generalize and the object part is too coarse to provide effective guidance.

**Quantitative Experiments.** We conduct an experiment to prove the proposed affordance are better than contact map in generalization, as shown in Tab. 8. The Point IOU metric is used for measure the consistency of prediction and the ground truth. The performance of contact maps drops markedly in testing, compared with affordance.

	<i>Open Set A</i>					<i>Open Set B</i>				
	<i>FID</i> ↓	<i>CD</i> ↓	<i>Top1</i> ↑	<i>Q1</i> ↑	<i>Pen.</i> ↓	<i>FID</i> ↓	<i>CD</i> ↓	<i>Top1</i> ↑	<i>Q1</i> ↑	<i>Pen.</i> ↓
DDIM [47]	0.247	4.13	0.455	0.0167	0.541	0.211	3.66	0.520	0.0167	0.542
DDPM [12]	0.254	3.97	0.461	0.0283	0.463	0.204	3.50	0.516	0.0128	0.602
Flow Matching	0.242	3.79	0.480	0.0240	0.501	0.176	2.76	0.538	0.0150	0.612

Table 6. The comparison of Flow Matching with Diffusion models with DDPM and DDIM.

	AFM	GFM	DDIM	DDPM
Time	$0.075 \pm 0.006$	$0.233 \pm 0.022$	$0.235 \pm 0.013$	$1.14 \pm 0.099$

Table 7. The inference time of Affordance Flow Matching (AFM) (first column); Grasp Flow Matching (second column); Grasp generation with DDIM and DDPM (The third and fourth columns).

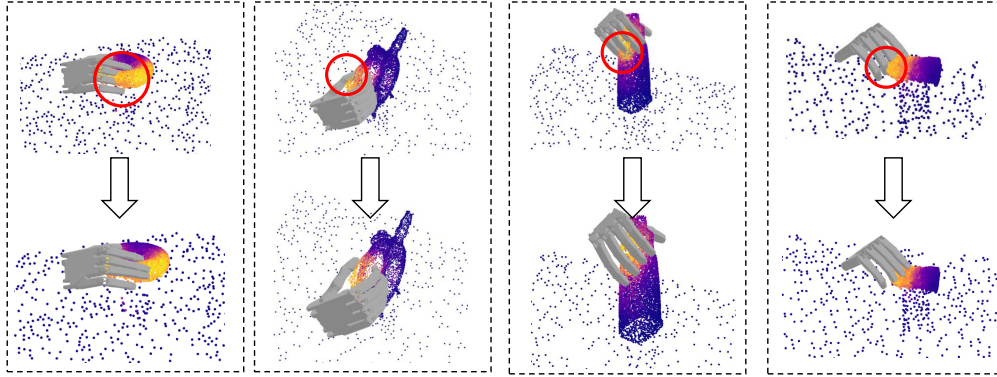


Figure 11. Visualization of grasps before and after affordance-guided grasp optimization.

### D.3. Qualitative Experiments of Affordance-guided Optimazation

To verify the effectiveness of the affordance-guided optimization, we offer more qualitative results. Figure 11 shows the grasping situations before and after optimization. It clearly shows that the QGC can prevent the grasp from penetrating the object and keep in line with the original intention.

### D.4. Extension to Dexterous Manipulation

We extend our methods to the manipulation task by cooperating the grasp ability of our framework with task planning ability of MLLM and the perception ability of computer vision foundation model to achieve dexterous manipulation task, like pouring water, following ReKep [15]. Specially, the MLLM first generates the relational keypoint constraints, then obtain an action solution by optimization solver. Then we employ our framework to generate a stable dexterous grasping and execute manipulation action according to the action solution. The visualization results can be found in Figure 14.

## E. Real World Experiments Details

**Experimental Environment** Figure 12 shows the settings of our real world experiments. The experiments are conducted on Leap hand, a Kinnova Gen3 6DOF arm and a Kinnova original wrist RGB-D camera. We collect several daily life objects to evaluate the open-set generalization in real world.

**Experiment Pipeline** To obtain the scene point clouds, we set the robotic arm to record partial point clouds along a specified trajectory, then fuse the local point clouds into a global point cloud and down-sample them to obtain the sampled point clouds, as shown in Figure 13. Next, the RGB image and user commands are fed into MLLM to obtain the guidance. The guidance and scene point clouds are then fed into our framework to sequentially generate affordance and grasp poses. During the grasp execution phase, we first move the robotic arm to the vicinity of the target, and then simultaneously control both the robotic arm and the dexterous hand to perform the grasp. The visualization can be found in Figure 15.



Point IOU $\uparrow$	Proposed Affordance		Contact Map	
	Train	Test	Train	Test
Set A	0.597	0.508	0.576	0.359
Set B	0.602	0.502	0.624	0.315

Table 8. Generation performance of the proposed affordance and contact map.

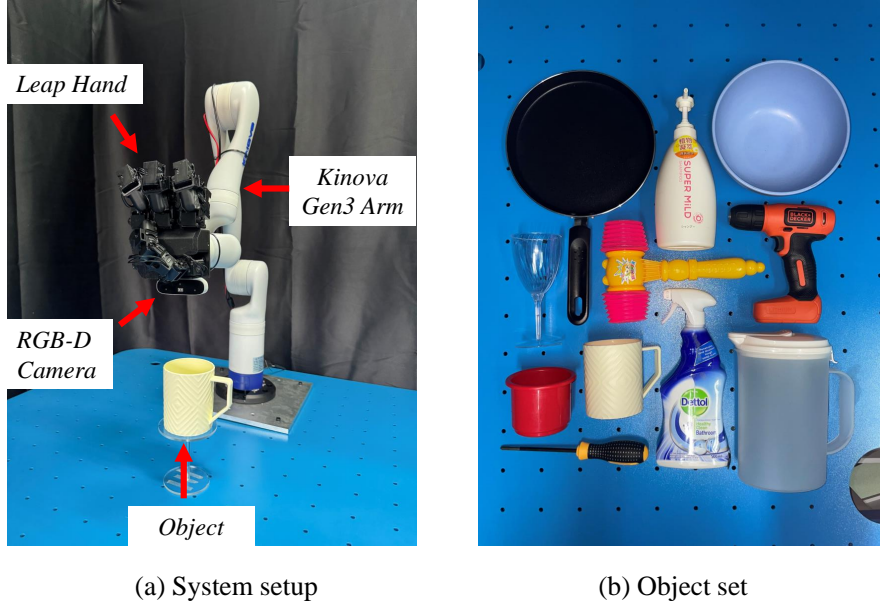


Figure 12. The illustration of our real world experiment settings.

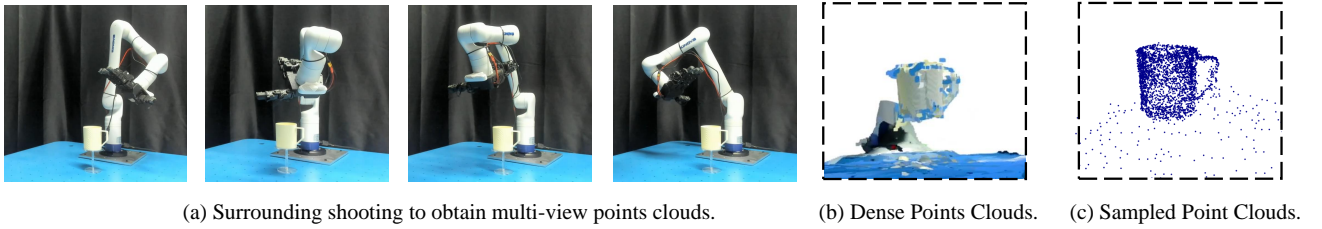


Figure 13. The visualization of the pipeline to obtain scene point clouds.

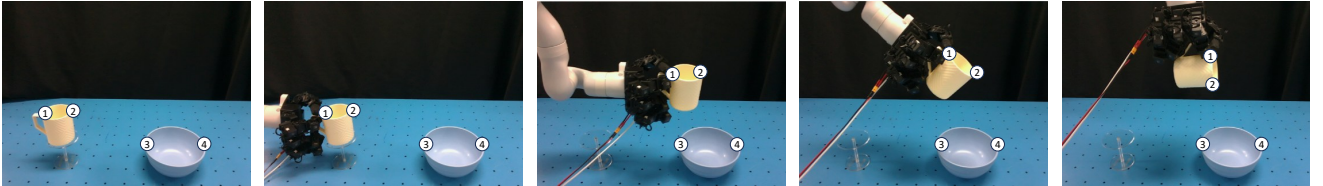


Figure 14. Visualization of the application of our method to the manipulation task.



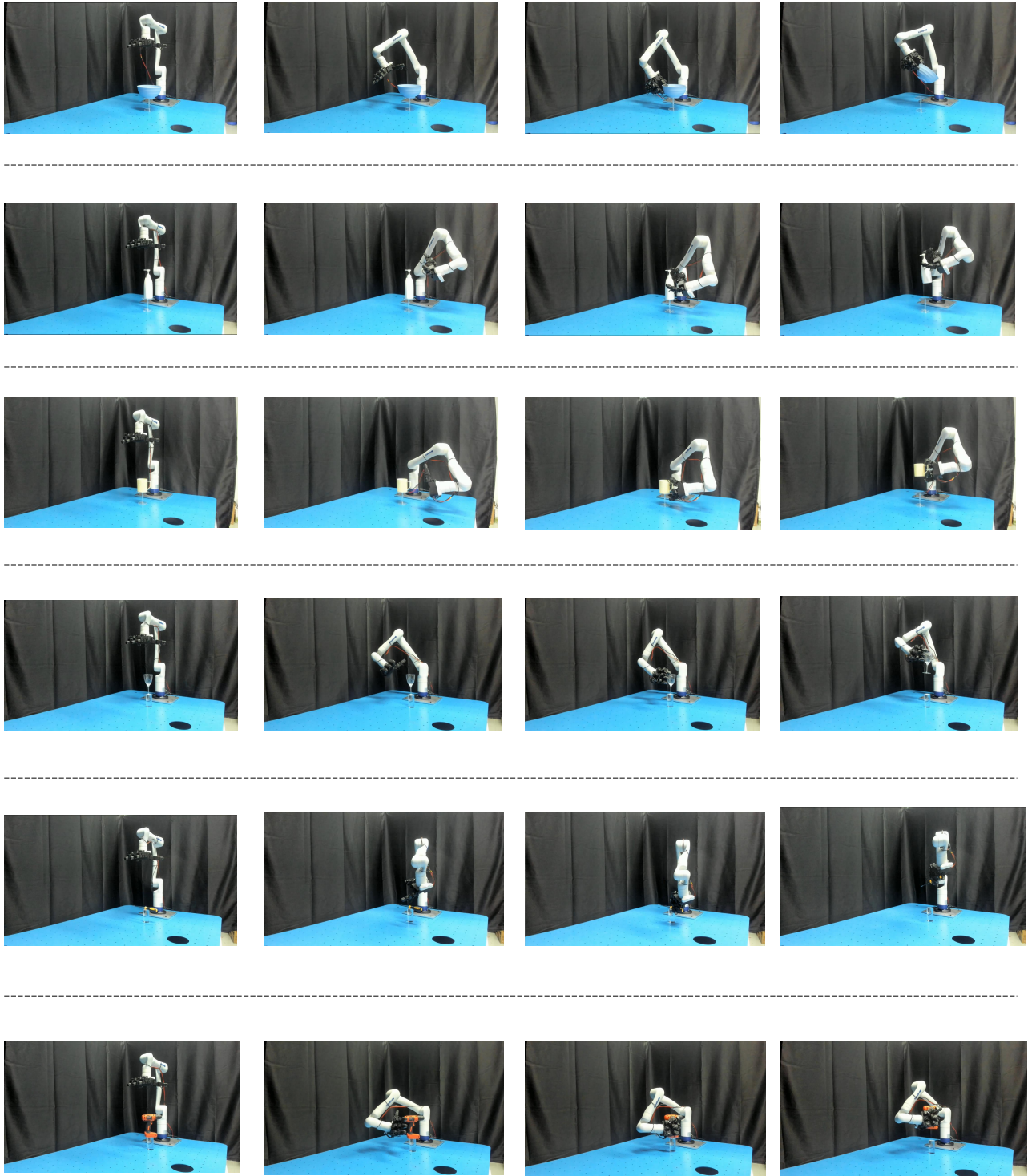


Figure 15. The visualization of motion in real world experiments.