

GEAL: Generalizable 3D Affordance Learning with Cross-Modal Consistency

Dongyue Lu Lingdong Kong Tianxin Huang Gim Hee Lee

National University of Singapore

dongyue.lu@u.nus.edu lingdong@comp.nus.edu.sg

{huangtx, gimhee.lee}@nus.edu.sg

<https://dylanorange.github.io/projects/geal>

Abstract

Identifying affordance regions on 3D objects from semantic cues is essential for robotics and human-machine interaction. However, existing 3D affordance learning methods struggle with generalization and robustness due to limited annotated data and a reliance on 3D backbones focused on geometric encoding, which often lack resilience to real-world noise and data corruption. We propose **GEAL**, a novel framework designed to enhance the generalization and robustness of 3D affordance learning by leveraging large-scale pre-trained 2D models. We employ a dual-branch architecture with Gaussian splatting to establish consistent mappings between 3D point clouds and 2D representations, enabling realistic 2D renderings from sparse point clouds. A granularity-adaptive fusion module and a 2D-3D consistency alignment module further strengthen cross-modal alignment and knowledge transfer, allowing the 3D branch to benefit from the rich semantics and generalization capacity of 2D models. To holistically assess the robustness, we introduce two new corruption-based benchmarks: PIAD-C and LASO-C. Extensive experiments on public datasets and our benchmarks show that **GEAL** consistently outperforms existing methods across seen and novel object categories, as well as corrupted data, demonstrating robust and adaptable affordance prediction under diverse conditions.

1. Introduction

3D affordance learning involves identifying interactive regions on objects given semantic cues such as image or textual instruction [7, 10], which is a fundamental competency for intelligent systems [6, 17] to infer how an object can be used or manipulated [24, 37, 63]. This understanding is vital for applications in robotics and human-machine interaction such as action prediction, object manipulation, and autonomous decision-making [5, 12, 13, 15]. For example, a robot equipped with affordance knowledge can intel-

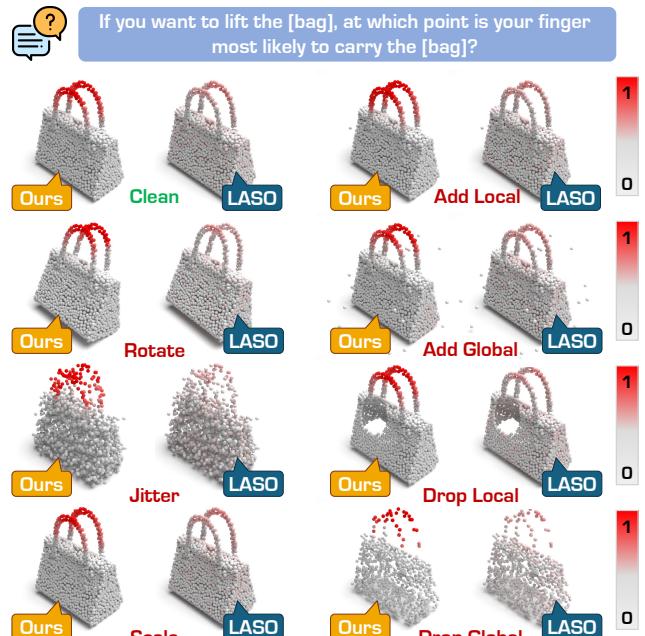


Figure 1. **3D affordance prediction under varied data noises.** Given a textual prompt, previous methods like LASO [28] (right side of each example) exhibit reduced robustness across different corruption types. In contrast, our proposed method, **GEAL** (left side of each example), maintains high accuracy and generalization across these challenging scenarios by effectively transferring knowledge from a large-scale pre-trained 2D foundation model, enhancing robustness and adaptability under diverse conditions.

ligently interacts with objects in its environment by determining where to grasp a handle or press a button.

Despite its potential, 3D affordance learning still faces significant challenges. Due to limited annotated data, 3D affordance models generally show poorer generalization than their 2D counterparts which benefit from abundant labeled data and large-scale pretraining [25]. Additionally, 3D models often rely on backbones that focus on positional and geometric encoding, limiting their capacity to capture global semantic content and making them vulner-

able to noisy or corrupted data from sensor inaccuracies, scene complexity, or processing artifacts in real-world settings [21, 50, 51]. These issues further hinder the robustness and adaptability of current 3D affordance learning methods.

In this paper, we introduce a novel framework **GEAL**, which is designed to enhance the generalization and robustness of 3D affordance learning through a dual-branch architecture that leverages the correspondence between 2D and 3D data. GEAL generates realistic 2D renderings directly from sparse 3D data by employing 3D Gaussian splatting (3DGS) [19] to build consistent mappings between 3D point clouds and 2D representations. This approach effectively creates a 2D branch from purely 3D data, which allows us to utilize the generalization capabilities and rich semantic knowledge of large-scale pre-trained 2D foundation models [46, 48] to enhance 3D affordance predictions.

We further introduce a granularity-adaptive fusion module, and a 2D-3D consistency alignment module to ensure robust multi-modal alignment. The granularity-adaptive fusion module dynamically integrates multi-level visual and textual features to address affordance queries at various scales and granularities. The 2D-3D consistency alignment module concurrently establishes reliable cross-modal correspondence with feature embeddings augmented to the Gaussian primitives of 3DGS, fostering effective knowledge transfer across branches, and enhances the generalization and robustness of the 3D branch by enforcing consistent alignment between 2D and 3D modalities.

In view of the limitation of data scarcity to benchmark the robustness of 3D affordance models, we create two datasets: **PIAD-Corrupt** and **LASO-Corrupt** from existing commonly used affordance datasets [28, 63]. We design these benchmark datasets by incorporating various types of real-world corruptions such as scaling, cropping, *etc.*, to ensure their suitability in evaluating the robustness of 3D affordance models. By contributing these benchmark datasets, we aim to fill a critical gap in the affordance learning community by providing a standard for evaluating the robustness of point cloud-based 3D affordance methods. Fig. 1 shows an example of the text description and the corresponding 3D affordance on 3D point clouds that are corrupted under various noise types.

We validate the generalization and robustness of our GEAL on both standard and corruption-based benchmarks, demonstrating that our approach consistently outperforms recent methods on all scenarios. Our experiments confirm that our GEAL effectively transfers knowledge from seen to unseen data and maintains high performance even under corruption, underscoring the adaptability of our framework across challenging scenarios.

Our main contributions are summarized as follows:

- We propose **GEAL**, a novel approach for generalizable 3D affordance learning. By employing 3DGS, we de-

velop a 2D affordance prediction branch for 3D point clouds, harnessing the robust generalization and semantic understanding of pre-trained 2D foundation models.

- We propose granularity-adaptive fusion and 2D-3D consistency alignment to integrate and propagate knowledge across the dual-branch architecture, and enhance the generalizability of the 3D branch using 2D knowledge.
- We establish two corruption-based benchmarks: **PIAD-C** and **LASO-C**, to holistically evaluate the robustness of 3D affordance learning under real-world scenarios, contributing a standard to the community for robustness analysis.
- Extensive experiments validate the strong performance of our approach on both mainstream and corruption 3D affordance learning benchmarks, proving its generalization ability and robustness across diverse conditions.

2. Related Work

2D Affordance Learning. Affordances refer to potential actions that objects or environments enable for an observer, based on their properties [5, 9, 22, 43]. Early methods for affordance detection mainly try to identify interaction regions in images and videos [26, 35, 52, 55, 64], though these often lacked precise localization of affordance-relevant object parts. To address this, later research improved affordance localization [4, 10, 24, 31, 34, 36, 37, 44, 63] given demonstration 2D data. Recently, large-scale pre-trained models [2, 48] have aligned visual features with affordance-related textual descriptions, reducing dependence on manual labels and enhancing affordance prediction in new contexts [34, 39, 40, 45]. Building on this, some studies [14, 24, 25] turn to leverage foundation models to generalize affordance detection to novel objects and views.

3D Affordance Learning. Extending affordance detection to 3D space presents challenges due to the need for accurate spatial and depth information. While some studies use 2D data to detect 3D affordance regions [5, 9, 26], they often struggle with precise 3D interaction sites. The availability of large-scale 3D object datasets [11, 30, 41] has driven efforts to map affordances directly onto 3D structures [8, 28, 42, 60, 63], aiming to capture complex spatial relationships. Recent methods [18, 28, 45] leverage 2D visual and language models for open-vocabulary affordance detection, enhancing generalization without fixed label sets. Despite these advancements, achieving robust generalization in 3D remains challenging, as 3D backbones still lack the generalization capabilities of 2D foundation models, thus, our method leverages large-scale 2D foundation models to improve 3D affordance generalization.

Robustness for 3D Affordance Learning. Real-world 3D affordance learning faces inevitable challenges from point cloud corruptions caused by scene complexity, sensor inaccuracies, and processing errors [16, 51]. Existing studies aim to improve and benchmark robustness against noise

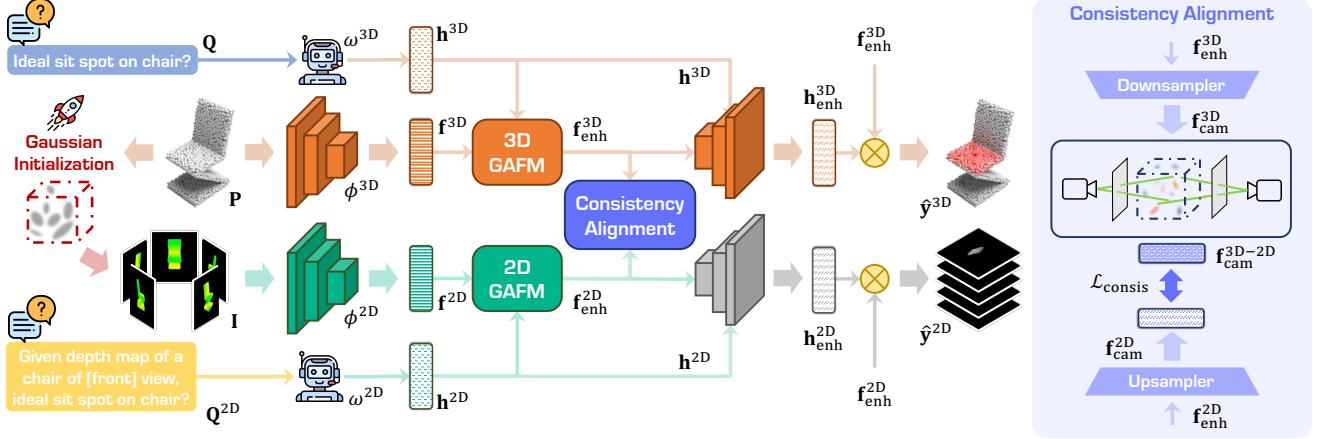


Figure 2. (Left): Framework Overview. The proposed **GEAL** consists of two branches: 3D and 2D. The 2D branch is established through 3D Gaussian Splatting to leverage the generalization capabilities of large pre-trained 2D models (*cf.* 3.1). We then perform cross-modality alignment, including **Granularity-Adaptive Visual-Textual Fusion** and **2D-3D Consistency Alignment**, to unify features from different modalities into a shared embedding space (*cf.* 3.2). Finally, we decode generalizable affordance from this embedding space (*cf.* 3.3). **(Right): Architecture of the 2D-3D Consistency Alignment Module.** This module maps features from 2D and 3D modalities into a shared embedding space and enforces consistency alignment to enable effective knowledge transfer across branches.

and corruption in 3D perception across real-world scenarios [20, 21, 23, 57, 58]. However, affordance learning uniquely requires precise identification of interactive regions under variable and degraded data conditions. To our knowledge, this work is the first to specifically address robustness in 3D affordance learning, providing a targeted approach to enhance reliability across diverse environments.

3. Methodology

In this section, we describe the technical components of our proposed **GEAL** framework. An overview of the full framework is shown in Fig. 2. Given an instruction Q and an object point cloud $\mathbf{P} \in \mathbb{R}^{N \times 3}$ with N points, GEAL predicts an affordance score $\mathbf{y} \in \mathbb{R}^N$, where each value in \mathbf{y} indicates the likelihood that a corresponding point supports the specified functionality. In Sec. 3.1, we employ Gaussian splatting as a cross-modal mapping to bridge the 2D and 3D modalities, establishing a 2D branch to leverage the generalization and robustness of large pre-trained 2D models. Sec. 3.2 details our cross-modal alignment strategy, incorporating both granularity-adaptive visual-text fusion and 2D-3D consistency alignment to unify these modalities in the embedding space. In Sec. 3.3, we outline the decoding process that derives robust and generalizable affordance predictions from the aligned feature space.

3.1. 3D-2D Mapping with Gaussian Splatting

Motivation. Current 3D affordance learning methods suffer from poor generalization due to limited annotated data and exhibit relatively weak robustness owing to limited global semantic capture. In contrast, 2D affordance learning methods [24, 25] leverage 2D foundation models [46, 48] pre-trained on large amounts of data, offering superior gen-

eralization and robustness. A 3D-2D mapping to leverage the strengths of 2D foundation models is thus promising. However, a direct projection of 3D point clouds onto 2D planes yields sparse 2D points without semantic and depth information that are not useful for feature extraction with 2D foundation models. To overcome this issue, we adopt 3D Gaussian Splatting [19] which represents 3D scenes as learnable Gaussian primitives for realistic, differentiable and high-speed rendering from arbitrary viewpoints. This approach allows us to synthesize realistic 2D images from sparse point clouds, preserving crucial semantic and depth information for downstream affordance learning tasks. Moreover, 3D Gaussian Splatting offers smoother transitions between points, preserves occlusions and depth cues for a coherent and accurate scene representation.

Gaussian Initialization. In 3D Gaussian Splatting, each Gaussian primitive is characterized by its 3D position μ represented by a 3D coordinate, a covariance matrix Σ which defines its shape and spread, spherical harmonic parameters c representing its color, and an opacity value α that indicates its transparency. To render 3D Gaussian primitives into 2D image planes, we apply point-based α -blending using a tile-based rasterizer for efficient rendering. The rendered color at each pixel v is calculated as follows:

$$C(v) = \sum_{i \in \mathcal{N}} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (1)$$

where c_i is the color of the i -th Gaussian, \mathcal{N} represents the Gaussians within the tile, and $\alpha_i = o_i G_i^{2D}(v)$. o_i is opacity of the i -th Gaussian and $G_i^{2D}(\cdot)$ represents the function of the i -th Gaussian projected onto 2D. Similarly, a depth map

can be rendered as:

$$D(v) = \sum_{i \in \mathcal{N}} d_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (2)$$

where d_i denotes the depth of the i -th Gaussian primitive under the provided camera pose.

To ensure that the rendered images accurately reflect the geometry of the input point cloud \mathbf{P} , we set the Gaussian mean positions to match the point coordinates, *i.e.* $\mu = \mathbf{P}$. The covariance Σ and opacity α are manually adjusted and then kept fixed during training to preserve the original geometry. Using the depth map from Eq. (2) with V camera poses and a predefined color map, we obtain realistic images $\mathbf{I} \in \mathbb{R}^{V \times 3 \times H \times W}$ that preserve both semantics and spatial information of the original point cloud, effectively bridging the 3D-2D gap. Treating the affordance score $y \in [0, 1]$ as grayscale color, we assign the color of each Gaussian to match its affordance score, *i.e.* $\mathbf{c} = \mathbf{y}$. We generate 2D affordance masks $\mathbf{y}_{2D} \in \mathbb{R}^{V \times H \times W}$, where each pixel represents the affordance score of the associated 3D point. This process establishes a coherent mapping from 3D point clouds and affordance scores to their 2D counterparts, using Gaussian splatting to generate realistic, informative 2D representations that enhance affordance learning.

Encoding. Our GEAL framework as shown in Fig. 2 comprises a 2D and a 3D branch with backbones $\phi^{2D}(\cdot)$ and $\phi^{3D}(\cdot)$, respectively. The 3D branch uses PointNet++ [47] for point cloud feature extraction, while the 2D branch employs DINOv2 [46] for image features. Both networks produce multi-scale features at various granularities. At each scale i , we extract features:

$$\mathbf{f}_i^{3D} = \phi_i^{3D}(\mathbf{P}), \quad \mathbf{f}_i^{2D} = \phi_i^{2D}(\mathbf{I}), \quad (3)$$

where $\mathbf{f}_i^{3D} \in \mathbb{R}^{B \times C_i^{3D} \times N_i^P}$ and $\mathbf{f}_i^{2D} \in \mathbb{R}^{B \times V \times C_i^{2D} \times N^I}$. B is the batch size, V is the number of views, and C_i^{3D} and C_i^{2D} are feature dimensions. N_i^P is the number of point in scale i , and N^I is image patch length. Note that the 3D spatial resolution N_i^P and C_i^{3D} differ between different scales due to the usage of PointNet++ [47].

We process the input prompt Q using lightweight language models $\omega^{3D}(\cdot)$ and $\omega^{2D}(\cdot)$ that share the same architecture. For the 2D input, we modify the prompt to Q^{2D} by adding: “**Given a depth map of a [object] in [view]**”, constructing **view-dependent prompt** to enhance context understanding. The text embeddings are:

$$\mathbf{h}^{3D} = \omega^{3D}(Q), \quad \mathbf{h}^{2D} = \omega^{2D}(Q^{2D}), \quad (4)$$

where $\mathbf{h}^{3D} \in \mathbb{R}^{B \times C^{txt} \times L}$ and $\mathbf{h}^{2D} \in \mathbb{R}^{B \times V \times C^{txt} \times L}$, with L as the sequence length.

3.2. Cross-Modal Consistency Alignment

Since point cloud, image, and text features are embedded in distinct spaces, we design alignment modules to map these

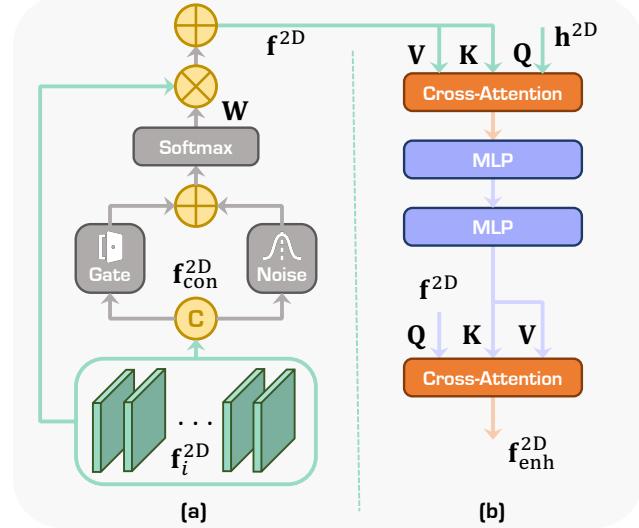


Figure 3. Illustration of the **Granularity-Adaptive Fusion Module**, it consists of a Flexible Granularity Feature Aggregation mechanism (a) and a Text-Conditioned Visual Alignment mechanism (b), we take the 2D branch as example.

multi-modal features into a shared embedding space. First, we fuse visual features at varying granularities with textual features through Granularity-Adaptive Visual-Textual Fusion, supporting affordance learning conditioned on instructions across different scales. Subsequently, we propagate knowledge from the 2D to 3D branch via a 2D-3D Consistency Alignment by enforcing consistency between 2D and 3D features.

Granularity-Adaptive Visual-Textual Fusion. Both 2D and 3D backbones capture different levels of granularity, with lower layers focusing on fine details and higher layers providing broader context. Since affordances can span multiple object parts, leveraging features at various granularities is advantageous. To achieve this, we introduce **Granularity-Adaptive Fusion Module (GAFM)**, which integrates multi-granularity visual features with textual cues via *Flexible Granularity Feature Aggregation* and *Text-Conditioned Visual Alignment*. These mechanisms allow the model to adaptively fuse features across different granularities, enhancing affordance prediction in response to specific instructions. An illustration of the Granularity-Adaptive Fusion Module is provided in Fig. 3.

Flexible Granularity Feature Aggregation. This mechanism aims to fuse visual features from different granularities. Taking the 2D branch as an example, we concatenate feature maps from the last m levels, forming an input tensor $\mathbf{f}_{con}^{2D} \in \mathbb{R}^{B \times V \times (m \times C^{2D}) \times N^I}$. Inspired by previous works [25, 65], we then compute adaptive soft weights to regulate the contribution of each feature level, enabling the model to adapt to varying levels of detail. These weights are computed via a gating function with learned noise, in-

troducing perturbations that enhance adaptability:

$$\mathbf{W} = \text{Softmax}(\mathbf{f}_{\text{con}}^{2D} \cdot \mathbf{W}_g + \sigma \cdot \epsilon), \quad (5)$$

where $\mathbf{W}_g \in \mathbb{R}^{(m \times C^{2D}) \times m}$ is a trainable weight matrix, $\mathbf{W} \in \mathbb{R}^{B \times m}$ represents the concatenation of weights w_i for each feature level, σ is a learned standard deviation controlling the noise scale, and $\epsilon \sim \mathcal{N}(0, 1)$ is Gaussian noise. These weights balance the influence of each feature level, enabling affordance reasoning across different granularities.

The fused feature map is then obtained by applying the adaptive weights to features from each level:

$$\mathbf{f}^{2D} = \sum_{i=1}^m w_i \odot \mathbf{f}_i^{2D}, \quad (6)$$

where \odot denotes element-wise multiplication and $w_i \in \mathbf{W}$. This adaptive aggregation yields a robust feature representation across varying conditions to enhance the generalization ability of the model.

Text-Conditioned Visual Alignment. This module is proposed to integrate visual features with the textual instruction. we follow [28, 56] to feed \mathbf{f}^{2D} and \mathbf{h}^{2D} into a transformer block. We first enhance the textual features \mathbf{h}^{2D} with visual features \mathbf{f}^{2D} through cross-attention, followed by refinement with two multilayer perceptrons (MLPs). We then acquire the visual features $\mathbf{f}_{\text{enh}}^{2D} \in \mathbb{R}^{B \times C^{\text{txt}} \times N^I}$ by querying the refined textual features with cross-attention. We thus ensure that the 2D visual features maintain their spatial structure while embedding the question-relevant information.

In the 3D branch, we align textual features with multi-granularity visual features in a similar manner. However, due to varying spatial resolutions and feature dimensions across scales in PointNet++ [47], directly concatenating all scales' features and computing soft weights as in Eq. (5) is not feasible. To address this, we first apply Text-Conditioned Visual Alignment to the 3D visual features at each scale, then upsample them to a uniform resolution. Finally, we perform Flexible Granularity Feature Aggregation on these upsampled features to produce the aggregated visual representation.

2D-3D Consistency Alignment. 2D features retain rich semantic context and strong generalization via the pre-trained backbone [46], while 3D features preserve geometric and spatial details, compensating for the loss of 2D information caused by self-occlusions. To propagate the knowledge inherently, we introduce **Consistency Alignment Module (CAM)** to ensure mutual alignment and knowledge transfer from 2D to 3D spaces.

Specifically, as shown in right part of Fig. 2, we map $\mathbf{f}_{\text{enh}}^{3D}$ and $\mathbf{f}_{\text{enh}}^{2D}$ into a shared embedding space. Given the 2D-3D correspondence, regions in 2D and 3D representations that map to the same spatial areas should exhibit similar feature representations. By enforcing this consistency, we facilitate

2D-3D knowledge propagation to enhance the understanding of affordances across both modalities of the model.

To align 3D and 2D features in the same embedding space, we employ a down-sampler consisting of two Conv1D layers that reduces the feature dimension of $\mathbf{f}_{\text{enh}}^{3D}$ to $\mathbf{f}_{\text{cam}}^{3D} \in \mathbb{R}^{B \times C^{\text{cam}} \times N}$. This processed feature acts as the representation for each point. We then leverage the established 2D-3D mapping using Gaussian splatting to project these point features into 2D. For each Gaussian, we treat its point feature vector as an inherent attribute. The 2D feature at pixel v is then rendered as:

$$\mathbf{F}(v) = \sum_{i \in \mathcal{N}} \mathbf{f}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (7)$$

where \mathbf{f}_i is the feature of the i -th Gaussian, α_i is its opacity, and $\mathbf{F}(v)$ is the resulting semantic feature at pixel v .

Similarly, we map the 2D features into the same embedding space using an up-sampler consisting of three Conv2D layers, which upsamples the spatial resolution of $\mathbf{f}_{\text{enh}}^{2D}$ while also reducing its feature dimension to $\mathbf{f}_{\text{cam}}^{2D} \in \mathbb{R}^{B \times V \times C^{\text{cam}} \times H \times W}$. V is the number of views and H and W is the spatial dimensions. Given the pixel positions from all V number of $H \times M$ feature maps as M , we can define the 3D-2D projected feature as $\mathbf{f}_{\text{cam}}^{3D-2D} = \{\mathbf{F}(v) | v \in M\}$. We then enforce a consistency constraint by minimizing the difference between the aligned 3D-2D features using L_2 loss:

$$\mathcal{L}_{\text{consis}} = \text{MSE}(\mathbf{f}_{\text{cam}}^{3D-2D}, \mathbf{f}_{\text{cam}}^{2D}). \quad (8)$$

This consistency loss $\mathcal{L}_{\text{consis}}$ encourages the model to maintain similar representations in both 2D and 3D spaces, effectively aligning affordance knowledge across domains. This alignment supports 2D-3D knowledge propagation, ensuring that the information learned in the 2D domain benefits the 3D features.

3.3. Decoding Generalizable Affordance

The affordance scores is decoded under the condition of affordance instructions. Through transformer decoder, the textual features attend to enhanced visual features, focusing the model on specific object parts for accurate predictions.

Our decoder architecture is shared across both 2D and 3D branches. In the 2D branch, textual features \mathbf{h}^{2D} and enhanced visual features $\mathbf{f}_{\text{enh}}^{2D}$ are processed through a 3-layer transformer decoder. Here, \mathbf{h}^{2D} serves as the query, and $\mathbf{f}_{\text{enh}}^{2D}$ acts as key and value, yielding updated textual features $\mathbf{h}_{\text{enh}}^{2D}$. Each layer comprises self-attention to refine textual relationships and cross-attention to guide focus toward relevant visual regions.

After the transformer decoder, these enhanced textual features serve as dynamic kernels to predict affordance scores from visual features. The final affordance prediction

\hat{y}^{2D} is obtained by multiplying h_{enh}^{2D} with f_{enh}^{2D} , followed by a sigmoid activation:

$$\hat{y}^{2D} = \text{sigmoid}(h_{enh}^{2D} \cdot f_{enh}^{2D}), \quad (9)$$

where $\hat{y}^{2D} \in \mathbb{R}^N$ denotes affordance scores.

Training. We employ a combination of Binary Cross-Entropy (BCE) and Dice loss to guide the affordance score prediction in each branch, addressing both class imbalance and segmentation accuracy. For the 2D branch, the loss function is:

$$\mathcal{L}^{2D} = \mathcal{L}_{BCE}^{2D} + \mathcal{L}_{Dice}^{2D}, \quad (10)$$

where \mathcal{L}_{BCE}^{2D} minimizes discrepancies between predicted and true affordance scores, and \mathcal{L}_{Dice}^{2D} improves the overlap between predicted and ground truth regions by maximizing intersection over union.

We adopt a two-stage training approach. We train the 2D branch in the first stage, optimizing it for robust feature extraction and affordance decoding. Except for the CAM (Consistency Alignment Module), all layers in the 2D branch are frozen in the second stage training. This approach allows the 3D branch to leverage fixed 2D features while adapting to 3D-specific characteristics. Consequently, the loss function for the 3D branch becomes:

$$\mathcal{L}^{3D} = \mathcal{L}_{BCE}^{3D} + \mathcal{L}_{Dice}^{3D} + \mathcal{L}_{consis}. \quad (11)$$

During inference, only the 3D branch is used, ensuring efficient and lightweight affordance prediction.

3.4. Corrupt Data Benchmark

To facilitate robust 3D affordance learning across diverse real-world scenarios, we establish two 3D affordance robustness benchmarks: **PIAD-C** and **LASO-C** based on the test sets of the commonly used datasets **PIAD** and **LASO** following [50]. We apply seven types of corruptions –*Add Global, Add Local, Drop Global, Drop Local, Rotate, Scale, and Jitter*–each with five severity levels. This results in a total of 4,890 object-affordance pairings, comprising 17 affordance categories and 23 object categories with 2,047 distinct object shapes. More details are provided in the supplementary material.

4. Experiment

4.1. Experimental Settings

Implementation Details. Our model is implemented in PyTorch and trained using the Adam optimizer with an initial learning rate of 1×10^{-4} for 50 epochs on a single NVIDIA A5000 GPU (24GB memory) with a batch size of 12. A step learning rate scheduler aids convergence. The 2D backbone DINOv2 [46] remains frozen during training, while the language model RoBERTa [32] is fine-tuned.

Type	Method	aIoU \uparrow	AUC \uparrow	SIM \uparrow	MAE \downarrow
Seen	MBDF [54]	9.3	74.9	0.415	0.143
	PMF [66]	10.1	75.1	0.425	0.141
	FRCNN [62]	12.0	76.1	0.429	0.136
	ILN [3]	11.5	75.8	0.427	0.137
	PFusion [61]	12.3	77.5	0.432	0.135
	XMF [1]	12.9	78.2	0.441	0.127
	IAGNet [63]	20.5	84.9	0.545	0.098
	LASO [28]	19.7	84.2	0.590	0.096
GEAL (Ours)		22.5	85.0	0.600	0.092
Unseen	MBDF [54]	4.2	58.2	0.325	0.213
	PMF [66]	4.7	60.3	0.330	0.211
	FRCNN [62]	5.1	61.9	0.332	0.195
	ILN [3]	4.7	59.7	0.325	0.207
	PFusion [61]	5.3	61.9	0.330	0.193
	XMF [1]	5.7	62.6	0.342	0.188
	IAGNet [63]	8.0	71.8	0.352	0.127
	LASO [28]	8.0	69.2	0.386	0.118
GEAL (Ours)		8.7	72.5	0.390	0.102

Table 1. The overall results of all comparative methods on **PIAD** [63]. **Seen** and **Unseen** are two partitions of the dataset. AUC and aIoU are shown in percentage. The **best** and **2nd best** scores from each metric are highlighted in **bold** and underlined, respectively.

Type	Method	aIoU \uparrow	AUC \uparrow	SIM \uparrow	MAE \downarrow
Seen	ReferTrans [27]	13.7	79.8	0.497	0.124
	ReLA [29]	15.2	78.9	0.532	0.118
	3D-SPS [38]	11.4	76.2	0.433	0.138
	IAGNet [63]	17.8	82.3	0.561	0.109
	LASO [28]	<u>20.8</u>	87.3	<u>0.629</u>	<u>0.093</u>
	LASO* [28]	19.7	85.2	0.600	0.097
	GEAL (Ours)	22.0	<u>86.7</u>	0.634	0.092
Unseen	ReferTrans [27]	10.2	69.1	0.432	0.145
	ReLA [29]	10.7	69.7	0.429	0.144
	3D-SPS [38]	7.9	68.8	0.402	0.158
	IAGNet [63]	12.9	77.8	0.443	0.129
	LASO [28]	14.6	80.2	0.507	0.119
	LASO* [28]	<u>15.6</u>	<u>79.9</u>	<u>0.549</u>	<u>0.108</u>
	GEAL (Ours)	16.7	80.9	0.567	0.106

Table 2. The overall results of all comparative methods on the **LASO** dataset [28]. **Seen** and **Unseen** are two partitions of the dataset. Results marked with * denote our reproduced results, following the data split reported in LASO [28]. AUC and aIoU are shown in percentage. The **best** and **2nd best** scores from each metric are highlighted in **bold** and underlined, respectively.

Datasets. We conduct experiments on **LASO**[28] and **PIAD**[63], both providing paired affordance and point cloud data. **LASO** contains 19,751 point cloud-question pairs over 8,434 objects (23 classes, 17 affordance categories) with **Seen** and **Unseen** splits to test generalization to novel affordance-object pairs. **PIAD** comprises 7,012 point clouds of the same categories, but some objects are entirely unseen during training, challenging the model’s generalization to novel object. Since PIAD lacks language annotations, we reuse LASO’s by randomly assigning questions to each affordance-object pair.

Type	aIoU↑		AUC↑		SIM↑		MAE↓	
	LASO	GEAL	LASO	GEAL	LASO	GEAL	LASO	GEAL
Scale	17.6	19.7	82.1	82.5	0.554	0.562	0.100	0.097
Jitter	14.7	17.0	80.3	80.6	0.501	0.505	0.103	0.099
Rotate	16.7	19.0	82.2	82.4	0.542	0.550	0.101	0.097
Drop-L	10.6	12.4	77.0	77.2	0.470	0.474	0.112	0.111
Drop-G	18.7	21.1	83.1	83.7	0.545	0.559	0.097	0.094
Add-L	15.7	18.5	81.0	81.1	0.525	0.536	0.100	0.095
Add-G	13.4	16.1	76.9	77.4	0.506	0.513	0.101	0.098

Table 3. Comparison of different methods under various corruption settings on the proposed PIAD-C benchmark, evaluated on the **Seen** partition. **Drop-L** denotes local drop, and **Drop-G** denotes global drop; similarly, **Add-L** and **Add-G** refer to local and global addition, respectively. AUC and aIoU are reported as percentages. For each metric, the **best** scores are highlighted in **bold**.

Type	aIoU↑		AUC↑		SIM↑		MAE↓	
	LASO	GEAL	LASO	GEAL	LASO	GEAL	LASO	GEAL
Scale	18.7	21.0	84.6	85.3	0.590	0.600	0.103	0.100
Jitter	15.4	17.8	81.3	81.9	0.516	0.517	0.107	0.106
Rotate	17.8	19.8	83.6	84.3	0.572	0.573	0.101	0.100
Drop-L	12.6	13.3	79.3	80.0	0.466	0.484	0.122	0.110
Drop-G	18.4	20.9	83.5	85.2	0.565	0.567	0.099	0.095
Add-L	17.6	20.2	82.7	84.4	0.566	0.572	0.103	0.100
Add-G	16.7	19.0	81.1	83.4	0.549	0.575	0.108	0.097

Table 4. Comparison of different methods under various corruption settings on the proposed LASO-C benchmark, evaluated on the **Seen** partition. **Drop-L** denotes local drop, and **Drop-G** denotes global drop; similarly, **Add-L** and **Add-G** refer to local and global addition, respectively. AUC and aIoU are reported as percentages. For each metric, the **best** scores are highlighted in **bold**.

Metrics. We use four metrics to assess performance: **AUC** [33] measures the ability to rank points correctly; **aIoU** [49] quantifies the overlap between predictions and ground truth; **SIM** [53] assesses the similarity by summing minimum values at each point; and **MAE** [59] calculates the average absolute difference between predictions and ground truth.

Baselines. We primarily compare our method with state-of-the-art approaches LASO [28] and IAGNet [63]. On PIAD, we evaluate against IAGNet and several image-point cloud cross-modal baselines, retraining LASO for comparison. On LASO, we compare with the original LASO method and other methods utilizing vision-language models for cross-modal alignment. To adapt IAGNet to LASO, its image backbone is replaced with a language model [28], keeping the rest of the architecture intact. Since the *Unseen* data split of LASO is not publicly available, we reproduce it based on the descriptions in their paper and report our results accordingly. Further experimental details are provided in the supplementary material.

4.2. Comparisons to State-of-the-Art Methods

Seen Categories: In Tab. 1 and Tab. 2, we present the performance of our **GEAL** compared to state-of-the-art approaches on the PIAD and LASO datasets under the *Seen* category setting. On the PIAD dataset, GEAL achieves the highest scores across all evaluation metrics, surpassing the previous best method, IAGNet [63]. Similarly, on the LASO dataset, GEAL outperforms LASO [28] on the ma-

Type	2D	3D	CAM	GAFM	aIoU↑	AUC↑	SIM↑	MAE↓
Seen	✓	✗	✗	✗	19.2	80.5	0.567	0.101
	✗	✓	✗	✗	19.5	83.5	0.585	0.097
	✓	✓	✓	✗	22.0	<u>84.4</u>	<u>0.592</u>	<u>0.094</u>
Unseen	✓	✓	✓	✓	22.5	85.0	0.600	0.092
	✓	✗	✗	✗	8.5	70.8	0.357	0.112
	✗	✓	✗	✗	8.0	69.2	<u>0.386</u>	0.118
Unseen	✓	✓	✓	✗	8.6	71.2	0.371	<u>0.105</u>
	✓	✓	✓	✓	8.7	72.5	0.390	0.102

Table 5. Ablation study on the impact of different components in **GEAL** on the PIAD dataset [63]. **Seen** and **Unseen** are two partitions of the dataset. **2D** denotes the use of the 2D baseline with a weighted sum mapping back to 3D. **3D** represents the 3D baseline. **CAM** is the consistency alignment module. **GAFM** is the granularity-adaptive fusion module. AUC and aIoU are shown in percentage. The **best** and second best scores from each metric are highlighted in **bold** and underlined, respectively.

Type	r	V	prompt	aIoU↑	AUC↑	SIM↑	MAE↓
Seen	112	6	✗	20.2	83.5	0.566	0.112
	112	12	✗	21.4	83.8	0.578	0.105
	112	12	✓	22.5	85.0	0.600	0.092
	112	14	✓	22.5	85.2	0.599	0.092
	224	14	✓	22.9	86.1	0.603	0.089
Unseen	112	6	✗	7.0	70.7	0.355	0.108
	112	12	✗	7.5	71.9	0.390	0.106
	112	12	✓	8.7	72.5	0.390	0.102
	112	14	✓	8.9	72.8	0.391	0.098
	224	14	✓	9.2	73.0	0.394	0.095

Table 6. Ablation study on the configuration of Gaussian Splatting parameters in **GEAL** on the PIAD dataset [63]. *r* denotes the resolution, *V* is the number of views, and **prompt** indicates whether a view-dependent prompt is used. **Seen** and **Unseen** are two partitions of the dataset. AUC and aIoU are shown in percentage. The **best** and 2nd best scores from each metric are highlighted in **bold** and underlined, respectively.

jority of metrics. These results highlight the effectiveness of GEAL in leveraging the rich semantic understanding from pre-trained 2D models through Gaussian splatting. The granularity-adaptive fusion and 2D-3D consistency alignment modules enable multi-granularity feature fusion and efficient knowledge transfer between the 2D and 3D modalities, enhancing the ability of the model to accurately predict affordance regions on seen categories.

Unseen Categories: The *Unseen* category setting evaluates the generalization ability of the model to novel objects not encountered during training. On both the PIAD (Tab. 1) and LASO (Tab. 2) datasets, **GEAL** consistently outperforms all baselines across metrics. Although the absolute performance values are lower due to the increased difficulty of unseen categories, GEAL maintains a performance edge over the baselines. This demonstrates that GEAL effectively generalizes to unseen categories, a result attributed to the integration of the 2D branch with a pretrained foundation model backbone and the cross-modal consistency alignment between the 2D and 3D branches. Qualitative comparisons with LASO on PIAD are shown in Fig. 4.

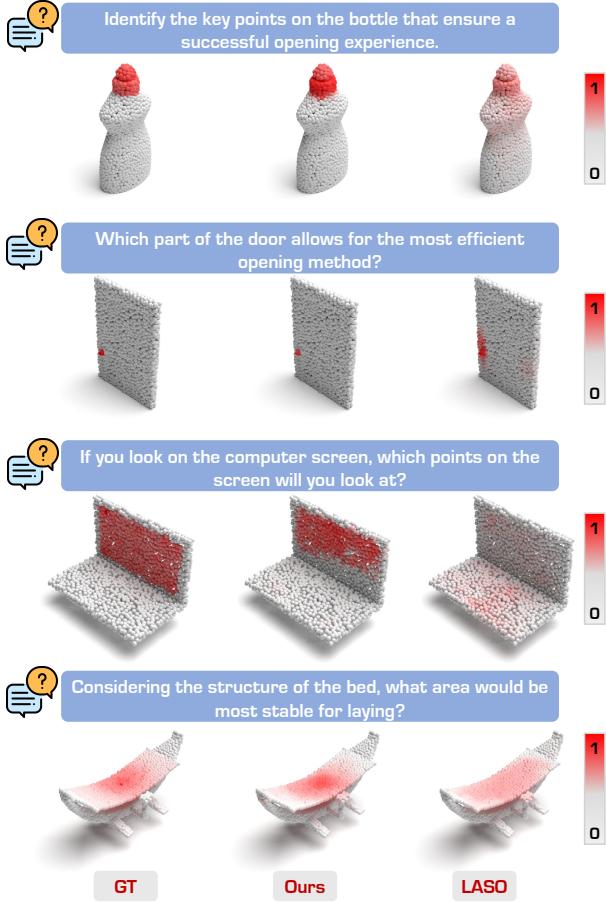


Figure 4. Qualitative comparisons between **GEAL** and **LASO** [28] on the PIAD [63] dataset. Top two rows displays results on *seen* partition, while bottom two rows show results on *unseen* partition. Our method demonstrates strong generalization on both seen and unseen partitions. *cf.* supplementary material for more examples.

Robustness on Corrupt Data: To assess robustness under real-world conditions, we compare **GEAL** with **LASO** on the proposed PIAD-C and LASO-C benchmarks after training on clean data. As shown in Tab. 3 and Tab. 4, **GEAL** consistently outperforms **LASO** across all corruption types and evaluation metrics. **GEAL** demonstrates superior resilience under various corruptions, achieving higher AUC and SIM scores while maintaining lower MAE values. This consistent outperformance indicates that the architecture of **GEAL** effectively mitigates the impact of data degradation. The robustness improvements are attributed to our dual-branch architecture and the 2D-3D consistency alignment module. By leveraging the robustness of pre-trained 2D models and enforcing cross-modal consistency, **GEAL** maintains high performance even when faced with corrupted or noisy 3D data.

4.3. Ablation Studies

Component Analysis. As shown in Table 5, we conduct an ablation study on the PIAD dataset [63] to evaluate the

effectiveness of each component in our proposed **GEAL** framework. We examine the impact of using only the 2D baseline with a weighted sum mapping back to 3D using the inverse process of Eq. (1)(2D), only the 3D baseline (3D), the consistency alignment module (CAM), and the granularity-adaptive fusion module (GAFM). Both the 2D and 3D baselines use only the last-layer features from their respective visual backbones to fuse with textual features without considering granularity. The results show that using only the 2D branch or only the 3D branch yields similar baseline performance. Integrating both branches with the consistency alignment module (CAM) leads to a noticeable improvement. Finally, our full model incorporating all components, including the granularity-adaptive fusion module (GAFM), achieves the best performance on both the *Seen* and *Unseen* sets. This demonstrates the effectiveness of our dual-branch architecture and underscores the importance of granularity-adaptive fusion and cross-modal consistency in enhancing the generalization capability of the model.

Gaussian Splatting Configuration Tuning. As shown in Table 6, we further investigate the impact of different Gaussian splatting configurations on the performance of our model. Specifically, we vary the rendering resolution (r), the number of views (V), and the inclusion of view-dependent prompts (**prompt**). The results indicate that increasing the number of views from 6 to 12 slightly improves performance, suggesting that additional viewpoints provide richer information for affordance learning. Incorporating view-dependent prompts significantly boosts performance, particularly on the *Seen* set, highlighting the importance of semantic guidance in our framework. Increasing the resolution from 112 to 224 yields only marginal gains, indicating that our model is robust to resolution changes and that higher resolutions offer diminishing returns. Balancing effectiveness and efficiency, we opt for a configuration of $r = 112$, $V = 12$, and the use of view-dependent prompts.

5. Conclusion

In conclusion, we present **GEAL**, a framework that improves the generalization and robustness of 3D affordance learning by leveraging large-scale pre-trained 2D models. Through a dual-branch architecture with Gaussian splatting, **GEAL** maps 3D point clouds to 2D representations, enabling realistic renderings from sparse data. The granularity-adaptive fusion and 2D-3D consistency alignment modules support cross-modal alignment and knowledge transfer, allowing the 3D branch to leverage rich semantic information from pre-trained 2D models. Experiments on public datasets and our corruption-based benchmarks show that **GEAL** consistently outperforms existing methods, demonstrating robust affordance prediction under varied conditions.

GEAL: Generalizable 3D Affordance Learning with Cross-Modal Consistency

Supplementary Material

Table of Contents

A Corrupt Data Benchmark	9
A.1. Corruption & Severity Level Settings	9
A.2 The PIAD-C Dataset	10
A.3 The LASO-C Dataset	10
B Benchmark Configuration	11
B.1. Datasets	11
B.2 Evaluation Metrics	13
B.3 Baselines	13
C Additional Quantitative Results	14
C.1. Complete Results on PIAD	14
C.2 Complete Results on LASO	14
D Additional Qualitative Results	14
D.1. Additional Qualitative Results on PIAD-C	14
D.2 Additional Qualitative Results on PIAD	14
E Broader Impact & Limitations	14
E.1. Societal Impact	14
E.2. Broader Impact	14
E.3. Potential Limitations	14
F Public Resource Used	15

A. Corrupt Data Benchmark

The robustness of models under real-world corruptions is a critical challenge in 3D point cloud analysis and 3D affordance learning [16, 20, 51, 57]. Unlike other 3D representations, point clouds often face various distortions caused by sensor inaccuracies, environmental complexities, and post-processing artifacts, which significantly impact downstream tasks [21, 23, 58]. For 3D affordance learning, ensuring robustness is paramount, as affordances are highly sensitive to object geometry and spatial details.

A.1. Corruption & Severity Level Settings

To standardize evaluation, we introduce a taxonomy of **seven atomic corruption types** – *Scale*, *Jitter*, *Rotate*, *Drop Global*, *Drop Local*, *Add Global*, *Add Local* – each simulating distinct real-world perturbations. These atomic corruptions simplify complex scenarios into controllable factors, enabling systematic analysis across **five levels of severity**. By providing a unified framework for benchmarking, we facilitate consistent and comprehensive assessment of model robustness, setting the stage for more resilient 3D affordance learning methods.

Below, we detail the construction methodology for each corruption type:

- **Jitter**

- *Description*: Adds Gaussian noise to perturb each point’s X, Y, and Z coordinates.
- *Mathematical Formulation*: For each point, a noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is added independently to X, Y, and Z.
- *Severity Levels*: The standard deviation σ varies as:

$$\sigma \in \{0.01, 0.02, 0.03, 0.04, 0.05\}$$

- **Scale**

- *Description*: Applies random scaling independently to the X, Y, and Z axes.
- *Mathematical Formulation*: Each axis is scaled by a factor $s \sim \mathcal{U}(\frac{1}{S}, S)$, where S determines the range of scaling.
- *Severity Levels*: The range of S is:

$$S \in \{1.6, 1.7, 1.8, 1.9, 2.0\}$$

After scaling, the point cloud is re-normalized to fit within a unit sphere.

- **Rotate**

- *Description*: Introduces random rotation to the point cloud.
- *Mathematical Formulation*: The rotation is specified by Euler angles (α, β, γ) , where:

$$\alpha, \beta, \gamma \sim \mathcal{U}(-\theta, \theta)$$

- *Severity Levels*: The angle range θ is:

$$\theta \in \{\pi/30, \pi/15, \pi/10, \pi/7.5, \pi/6\}$$

This approach does not guarantee uniform sampling in $\text{SO}(3)$, but provides sufficient variation to simulate diverse rotations.

- **Drop Global**

- *Description*: Randomly removes a percentage of points from the point cloud.
- *Method*: Shuffle all points and drop the last $N \cdot \rho$ points, where $N = 2048$ is the total number of points.
- *Severity Levels*: The proportion ρ is:

$$\rho \in \{0.25, 0.375, 0.5, 0.675, 0.75\}$$

- **Drop Local**

- *Description*: Removes points in clusters around randomly selected local regions.
- *Method*:

1. Randomly choose the number of local regions $C \sim \mathcal{U}\{1, 8\}$.
 2. For each region i :
 - * Randomly select a local center.
 - * Assign a cluster size N_i .
 - * Drop the N_i -nearest neighbor points to the center.
 3. Repeat for C regions.
- *Severity Levels*: The total number of points to drop K is:
- $$K \in \{100, 200, 300, 400, 500\}$$

• Add Global

- *Description*: Uniformly samples additional points inside a unit sphere and appends them to the point cloud. The added points are treated as noise and assigned a label of 0.
- *Method*: Sample K random points within a unit sphere.
- *Severity Levels*: The total number of added points K is:

$$K \in \{10, 20, 30, 40, 50\}$$

• Add Local

- *Description*: Adds clusters of points around randomly selected local regions. The added points are treated as noise and assigned a label of 0.
- *Method*:
 1. Shuffle points and select $C \sim \mathcal{U}\{1, 8\}$ as the number of local centers.
 2. For each center i :
 - * Define a cluster size N_i .
 - * Generate neighboring points' coordinates from:

$$\mathcal{N}(\mu_i, \sigma_i^2 I)$$

where μ_i is the i -th local center, and $\sigma_i \sim \mathcal{U}(0.075, 0.125)$.

3. Append generated points to the cloud one cluster at a time.
- *Severity Levels*: The total number of added points K is:

$$K \in \{100, 200, 300, 400, 500\}$$

A.2. The PIAD-C Dataset

Our proposed PIAD-C dataset is constructed from the test set of the **Seen** partition in PIAD [63], specifically designed to evaluate the robustness of affordance detection models under various corruption scenarios. This dataset includes a total of 2,474 object-affordance pairings, representing 17 affordance categories and 23 object categories, and with 1,012 distinct clean object shapes. Comprehensive statistics, detailing object categories, their corresponding affordance categories, and the number of object-affordance pairings, are presented in Tab. 7. We include additional visualization examples for the PIAD-C dataset in Fig. 5.

#	Object Category	Affordance	Number
1	Earphone •	listen, grasp	70
2	Bag •	contain, open, grasp, lift	50
3	Chair •	move, support, sit	587
4	Refrigerator •	contain, open	53
5	Knife •	stab, cut, grasp	138
6	Dishwasher •	contain, open	39
7	Keyboard •	press	25
8	Scissors •	stab, cut, grasp	29
9	Table •	move, support	194
10	StorageFurniture •	contain, open	92
11	Bottle •	contain, wrap_grasp, open, grasp, pour	273
12	Bowl •	contain, wrap_grasp, pour	83
13	Microwave •	contain, open	47
14	Display •	display	52
15	TrashCan •	contain, open, pour	69
16	Hat •	wear, grasp	66
17	Clock •	display	9
18	Door •	open, push	47
19	Mug •	contain, wrap_grasp, grasp, pour	126
20	Faucet •	open, grasp	95
21	Vase •	contain, wrap_grasp, pour	134
22	Laptop •	press, display	112
23	Bed •	lay, support, sit	84
Total	23	17	2474

Table 7. Detailed statistics of the **PIAD-C** dataset, showing the object categories, their corresponding affordance categories, and the number of object-affordance pairings for each category.

#	Object Category	Affordance	Number
1	Door •	open, push, pull	35
2	Clock •	display	34
3	Dishwasher •	open, contain	20
4	Earphone •	listen, grasp	28
5	Vase •	contain, pour, wrap_grasp	167
6	Knife •	stab, grasp, cut	59
7	Bowl •	contain, pour, wrap_grasp	36
8	Bag •	open, contain, lift, grasp	25
9	Faucet •	open, grasp	80
10	Scissors •	stab, grasp, cut	11
11	Display •	display	58
12	Chair •	sit, support, move	858
13	Bottle •	grasp, wrap_grasp, open, contain, pour	122
14	Microwave •	open, contain	23
15	StorageFurniture •	open, contain	183
16	Refrigerator •	open, contain	23
17	Mug •	contain, grasp, pour, wrap_grasp	45
18	Keyboard •	press	10
19	Table •	support, move	431
20	Bed •	sit, support, lay	36
21	Hat •	wear, grasp	26
22	Laptop •	display, press	55
23	TrashCan •	open, contain, pour	51
Total	23	17	2416

Table 8. Detailed statistics of the **LASO-C** dataset, showing the object categories, their corresponding affordance categories, and the number of distinct objects for each category.

A.3. The LASO-C Dataset

Our proposed LASO-C dataset is derived from the test set of the **Seen** partition in LASO [28], focusing on evaluating model robustness against point cloud corruptions. This dataset comprises 2,416 object-affordance pairings, covering 17 affordance categories and 23 object categories, with a total of 1,035 distinct clean object shapes. Comprehensive

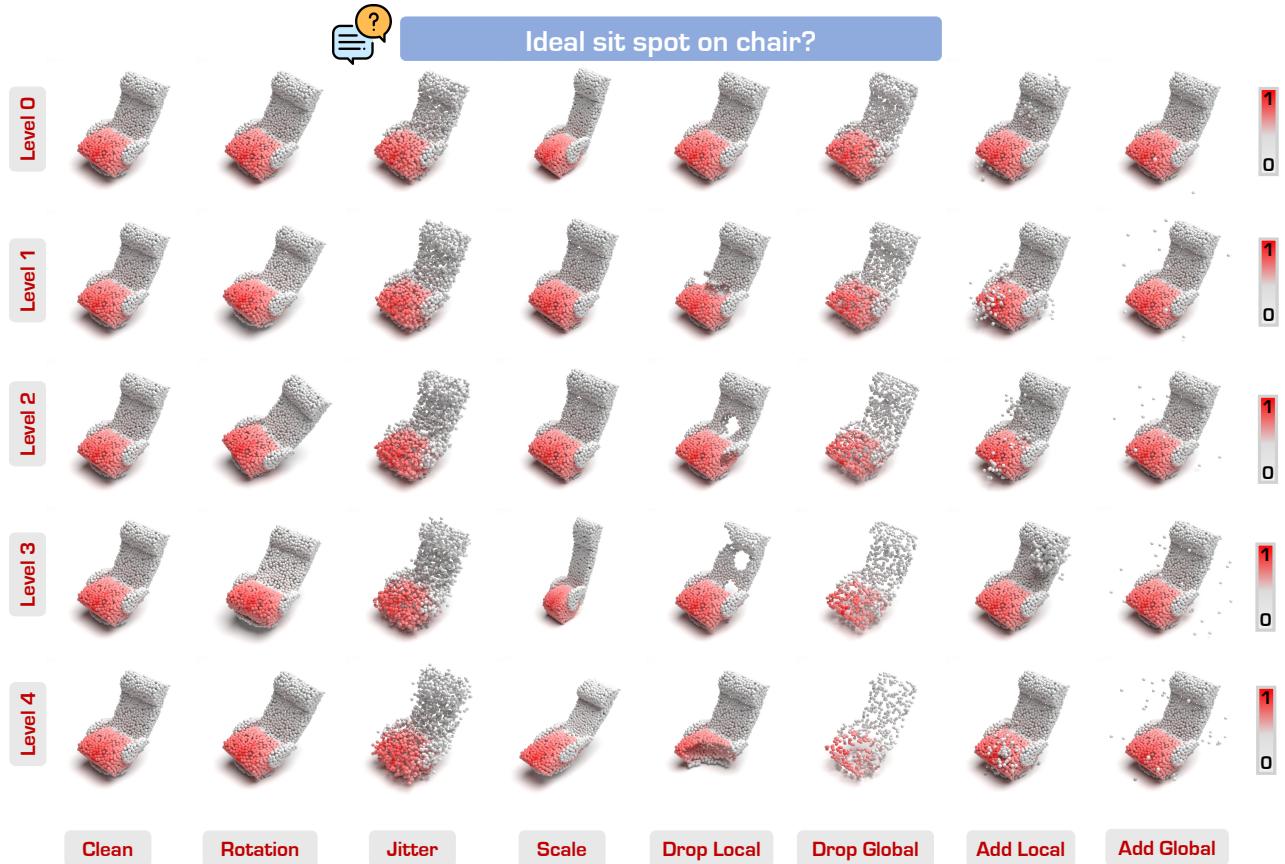


Figure 5. Visualization examples of the PIAD-C dataset. We show 7 corruption types across 5 severity levels.

statistics, detailing object categories, their corresponding affordance categories, and the number of object-affordance pairings, are presented in Tab. 8. We include additional visualization examples for the LASO-C dateset in Fig. 6.

B. Benchmark Configuration

B.1. Datasets

We conduct experiments primarily on the **LASO**[28] and **PIAD**[63] datasets, both of which provide paired affordance and point cloud data for evaluating 3D affordance learning.

LASO. This dataset is a pioneering benchmark designed to enable language-guided affordance segmentation of 3D objects. It includes **19,751 point cloud-question pairs** across **8,434 unique object shapes**, spanning **23 object categories** and **17 affordance types**. Derived from **3D-AffordanceNet** [8], the dataset pairs 3D object point clouds with questions that were carefully crafted by human experts and augmented using **GPT-4**. This process incorporates principles of *contextual enrichment*, *concise phrasing*, and *structural diversity*, enhancing the linguistic variety and

complexity of the dataset.

The LASO dataset introduces two distinct evaluation settings:

- **Seen Setting:** Models are trained and tested on overlapping object-affordance combinations, ensuring that both the object classes and affordance types in the training set are also present in the test set.
- **Unseen Setting:** This setting is designed to evaluate generalization capabilities. Certain object-affordance combinations (*e.g.*, “grasp-mug”) are excluded during training but appear in testing. This setting challenges models to transfer affordance knowledge learned from seen combinations (*e.g.*, “grasp-bag”) to novel combinations, promoting robust generalization.

These settings promote a comprehensive evaluation of models’ abilities to generalize affordance knowledge to unseen object-affordance pairings, a critical aspect for real-world deployment. The dataset also emphasizes diverse affordance scales and shapes, presenting significant challenges for perception models. By addressing the semantic limitations of traditional visual-only 3D affordance

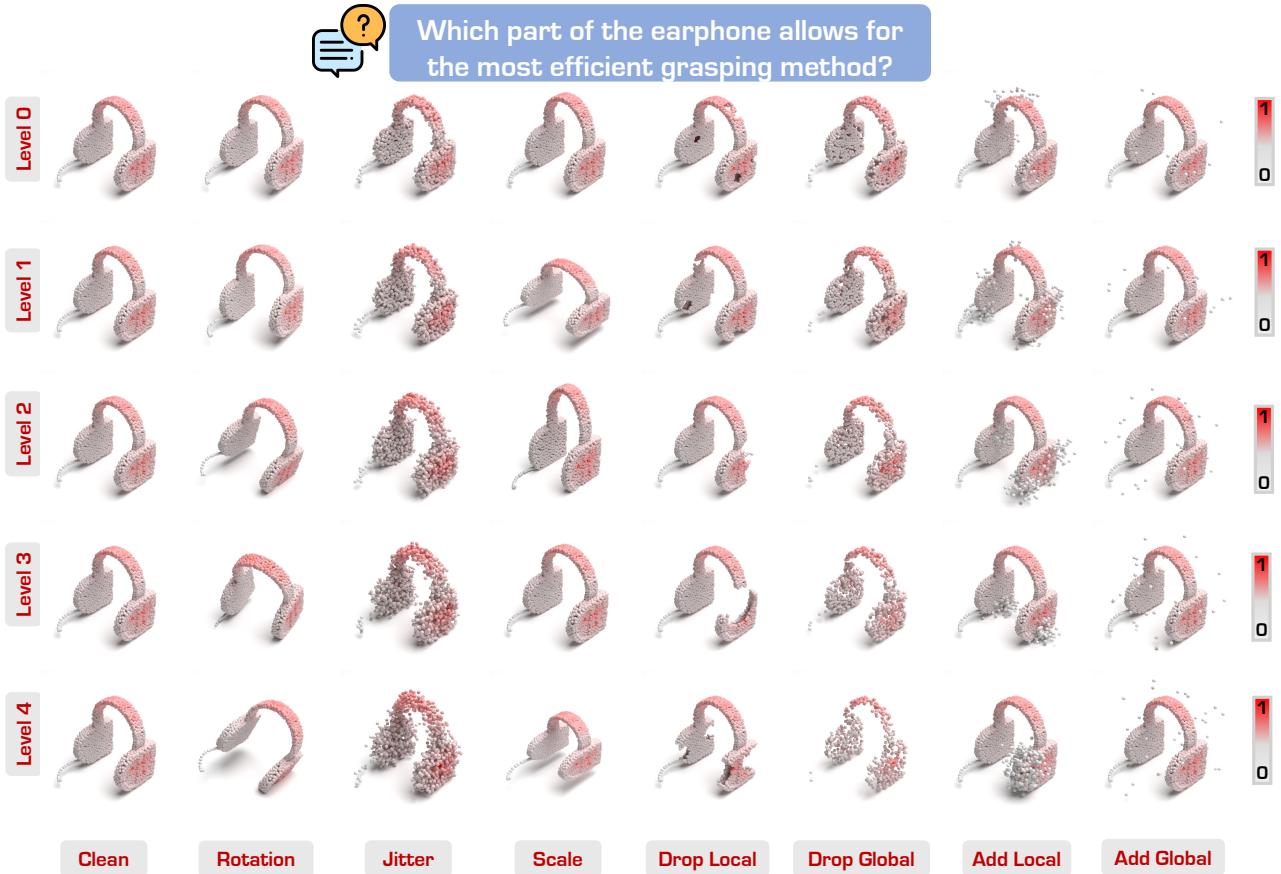


Figure 6. Visualization examples of the LASO-C dataset. We show 7 corruption types across 5 severity levels.

datasets, LASO bridges the gap between 3D perception and natural language understanding, encouraging cross-modal learning. This integration fosters advancements in embodied AI, enabling tasks that require nuanced reasoning and action in real-world environments.

PIAD. The Point-Image Affordance Dataset (PIAD) [63] is specifically curated to advance the task of grounding 3D object affordances using 2D interactions. PIAD consists of **7,012 point clouds** and **5,162 images**, spanning **23 object classes** and **17 affordance categories**. Unlike other datasets, PIAD pairs point clouds with images that demonstrate corresponding affordances. For example, a point cloud of a “Chair” affords “Sit,” and its paired image depicts a person sitting on a chair. These cross-modal pairings ensure consistency in affordance relationships while leveraging distinct modalities.

PIAD introduces two distinct evaluation settings:

- **Seen Setting:** In this setting, both objects and affordances in the training and testing sets are consistent. Point clouds and images of the same object categories and affordance types are included during training, allowing models to

learn affordance relationships in a supervised manner. This standard evaluation setting enables benchmarking on familiar object-affordance combinations.

- **Unseen Setting:** The Unseen partition presents a more challenging evaluation by excluding certain object categories from the training set entirely. For instance, some object categories are entirely unseen during training. This partition tests the ability of methods to transfer affordance knowledge across completely novel object instances and contexts, simulating real-world scenarios where interaction data is sparse or varied.

Annotations in PIAD include detailed affordance labels for point clouds, represented as heatmaps indicating the likelihood of affordance at each point. Paired images are annotated with bounding boxes for interactive subjects and objects, along with affordance category labels. This comprehensive annotation schema supports diverse affordance-learning paradigms and provides a robust benchmark for evaluating models in both Seen and Unseen scenarios.

Note that PIAD does not include language annotations. Since PIAD and LASO share the same object classes, affor-

dance categories, and same 58 affordance-object pairings, we reuse LASO’s language annotations for PIAD. For each object and affordance category label in PIAD, we randomly select a question from LASO’s question dataset corresponding to that affordance-object pairing.

B.2. Evaluation Metrics

To comprehensively evaluate the performance of our method, we employ four widely used metrics: **AUC**, **aIoU**, **SIM**, and **MAE**. Each metric is designed to assess different aspects of affordance prediction, providing a robust and multi-faceted evaluation framework. Below, we detail the formulation and significance of each metric:

- **Area Under the ROC Curve (AUC)** [33]: AUC measures the model’s ability to distinguish between regions of high and low affordance saliency on the point cloud. Specifically, the saliency map is treated as a binary classifier at various threshold levels, and a Receiver Operating Characteristic (ROC) curve is generated by plotting the true positive rate (TPR) against the false positive rate (FPR) at each threshold. AUC provides a single scalar value summarizing the overall performance, where higher values indicate better discrimination ability. It is particularly useful for comparing models’ effectiveness in highlighting affordance-relevant regions.
- **Average Intersection over Union (aIoU)** [49]: IoU is a standard metric for comparing the similarity between two arbitrary regions—in this case, the predicted affordance region and the ground truth. It is defined as the size of the intersection between the two regions divided by the size of their union:

$$IoU = \frac{TP}{TP + FP + FN}, \quad (12)$$

where TP , FP , and FN denote true positives, false positives, and false negatives, respectively. The aIoU extends this metric to compute the average IoU across all categories and test samples, providing a quantitative measure of the overlap between predicted and labeled affordance regions. Higher values indicate better alignment between the prediction and the ground truth.

- **Similarity (SIM)** [53]: The SIM metric evaluates how closely the predicted affordance map matches the ground truth by summing the minimum values at each point. For normalized prediction and ground truth maps P and Q , the similarity is calculated as:

$$SIM(P, Q) = \sum_i \min(P_i, Q_i), \quad (13)$$

where the inputs are normalized such that $\sum_i P_i = \sum_i Q_i = 1$. SIM provides a measure of how well the model captures the relative affordance distribution across the point cloud. A higher similarity score reflects greater

consistency between the predicted and true maps, making it a valuable metric for evaluating spatial prediction fidelity.

- **Mean Absolute Error (MAE)** [59]: MAE quantifies the average absolute difference between the predicted affordance values and the ground truth, offering a direct measure of prediction accuracy. For n points in a point cloud, it is calculated as:

$$MAE = \frac{1}{n} \sum_{i=1}^n |e_i|, \quad (14)$$

where e_i is the point-wise error. MAE is particularly useful for evaluating overall prediction quality by penalizing larger deviations. Lower MAE values indicate better performance, as they reflect a smaller error margin between the predicted and ground truth affordance scores.

Together, these metrics provide a comprehensive framework to benchmark the performance of affordance prediction models. AUC evaluates ranking capability, aIoU measures spatial overlap, SIM assesses prediction similarity, and MAE quantifies overall prediction accuracy. By combining these complementary metrics, we ensure a holistic evaluation of model performance under diverse scenarios.

B.3. Baselines

We evaluate our method against state-of-the-art approaches on both the PIAD and LASO datasets. Among these, LASO [28] is the closest to our method, as it also generates affordance scores based on textual cues. Additionally, we include 3D cross-modal baselines such as 3D-SPS [38], and image segmentation methods like ReferTrans [27] and RelA [29], which leverage vision-language models for cross-modal alignment. Results for these methods are referenced directly from the LASO paper.

On the PIAD dataset, we compare against IAGNet [63], a method that grounds 3D affordances by transferring knowledge from demonstration images into point clouds. Furthermore, this benchmark includes advanced image-point cloud cross-modal methods, including MBDF [54], PMF [66], FRCNN [62], ILN [3], PFusion [61], and XMF [1]. These baselines align image and point cloud features in various ways. Results for these baselines are taken from the IAGNet paper, except for LASO, which is retrained in the PIAD setting.

Below is a brief introduction to the baselines:

- **LASO** [28]: Generates affordance segmentation masks using textual-conditioned affordance queries, focusing on cross-modal alignment between text and 3D objects.
- **IAGNet** [63]: Grounds 3D affordances by transferring knowledge from 2D demonstration images to point clouds, leveraging cross-modal affordance reasoning.

- **3D-SPS** [38]: A 3D visual grounding method that selects linguistic keypoints for affordance segmentation, adapted by removing its bounding box prediction module.
- **ReLA** [29]: Originally designed for image-based referring expression segmentation, it segments point clouds based on language expressions by adapting image region features to grouped point features.
- **ReferTrans** [27]: A transformer-based architecture for image-based expression segmentation, modified for point clouds by replacing the image backbone with a 3D backbone and focusing solely on mask prediction.
- **MBDF-Net (MBDF)** [54]: Employs an Adaptive Attention Fusion (AAF) module for cross-modal feature fusion, with modifications to exclude camera intrinsic parameters.
- **PMF** [66]: Uses a residual-based fusion model to combine image and point cloud features, incorporating convolution and attention, while omitting perspective projection.
- **FusionRCNN (FRCNN)** [62]: Fuses proposals extracted from images and point clouds through iterative self-attention and cross-attention mechanisms.
- **ImloveNet (ILN)** [3]: Projects image features into 3D space using a learnable mapping, and fuses these with point cloud features using an attention mechanism.
- **PointFusion (PFusion)** [61]: Performs dense fusion by combining global and point-wise features extracted separately from point clouds and images.
- **XMFnet (XMF)** [1]: Fuses localized features from point clouds and images using a combination of cross-attention and self-attention, originally designed for cross-modal point cloud completion.

C. Additional Quantitative Results

C.1. Complete Results on PIAD

The complete results of the comparative methods for all object categories in the **Seen** and **Unseen** partitions of the PIAD dataset [63] are provided in Tab. 9 and Tab. 10, respectively.

C.2. Complete Results on LASO

The complete results of the comparative methods for all object categories in the **Seen** and **Unseen** partitions of the LASO dataset [28] are provided in Tab. 11 and Tab. 12, respectively.

D. Additional Qualitative Results

D.1. Additional Qualitative Results on PIAD-C

We include additional qualitative results of **GEAL** and **LASO** [28] on the PIAD-C dataset in Fig. 7.

D.2. Additional Qualitative Results on PIAD

We include additional qualitative results of **GEAL** and **LASO** [28] on the PIAD dataset in Fig. 8.

E. Broader Impact & Limitations

E.1. Societal Impact

The proposed framework for 3D affordance learning has significant societal implications, enabling embodied intelligence for effective robot and AI interaction with surroundings. This advancement can enhance automated systems' efficiency and safety in fields like healthcare, elderly care, and disaster response, where understanding object affordances is critical. This technology also has the potential to empower individuals with disabilities by enabling assistive robots to perform tasks such as fetching, opening, or manipulating objects. Applications in education and augmented or virtual reality could transform learning and entertainment by offering immersive and interactive experiences.

E.2. Broader Impact

Affordance learning can redefine robotics automation by improving autonomy and adaptability in industries. In manufacturing, it allows robots to handle diverse objects with minimal reprogramming, optimizing workflows and reducing human workload. In agriculture and environmental monitoring, affordance-aware systems can adapt to dynamic environments for precise operations. Integrating affordance grounding with augmented and virtual reality enables new possibilities in training, simulation, and interactive applications. This could drive innovations in medical training, such as AR-guided surgeries, and in gaming, offering intuitive and immersive user experiences through affordance-based interactions.

E.3. Potential Limitations

Despite its advantages, the proposed framework may encounter certain limitations:

- **Limited Generalization for Internal Affordances:** The framework struggles to accurately perceive and generalize affordances associated with the internal properties of objects, such as the "contain" affordance of a bottle. This limitation arises because point cloud processing primarily focuses on an object's external surface, often neglecting internal structures. Furthermore, the scarcity of high-quality data representing internal affordances, hampers the system's ability to generalize on such affordances.
- **Ethical Concerns:** In applications such as surveillance or autonomous decision-making, the deployment of the framework introduces potential ethical concerns. Misuse of the technology could infringe on privacy or lead to a lack of accountability in critical decision-making sce-

#	Category	LASO [28]				GEAL (Ours)				
		aIOU↑	AUC↑	SIM↑	MAE↓	aIOU↑	AUC↑	SIM↑	MAE↓	
1	StorageFurniture ●	Bag ●	23.4	83.3	0.567	0.090	24.0	85.1	0.588	0.088
2		Bed ●	21.1	87.3	0.587	0.097	22.7	88.1	0.595	0.091
3		Bowl ●	7.4	76.2	0.736	0.114	9.8	84.1	0.761	0.105
4		Clock ●	7.5	91.5	0.473	0.077	11.1	92.5	0.596	0.051
5		Dishwash ●	24.7	91.9	0.464	0.069	26.2	92.9	0.496	0.058
6		Display ●	32.5	91.5	0.719	0.083	37.7	91.3	0.726	0.104
7		Door ●	10.1	81.2	0.437	0.064	11.0	83.8	0.395	0.054
8		Earphone ●	18.8	85.9	0.615	0.094	21.6	87.6	0.654	0.086
9		Faucet ●	19.9	79.9	0.517	0.099	19.1	83.6	0.602	0.078
10		Hat ●	4.7	65.9	0.604	0.148	7.8	74.2	0.620	0.133
11		Keyboard ●	17.3	87.2	0.419	0.077	20.8	87.5	0.430	0.065
12		Knife ●	14.8	81.2	0.249	0.059	15.2	84.6	0.257	0.048
13		Laptop ●	15.5	89.8	0.671	0.060	23.5	94.1	0.717	0.046
14		Microwave ●	29.2	94.1	0.566	0.072	31.2	94.2	0.575	0.069
15		Mug ●	30.1	96.8	0.524	0.037	35.5	96.9	0.545	0.037
16		Refrigerator ●	10.7	76.5	0.578	0.107	17.5	77.2	0.607	0.091
17		Scissors ●	23.2	87.1	0.473	0.070	24.7	89.6	0.460	0.070
18		Chair ●	27.5	88.1	0.649	0.094	28.5	89.0	0.652	0.066
19		Table ●	24.1	91.2	0.631	0.055	31.9	95.8	0.698	0.040
20		TrashCan ●	10.1	78.2	0.627	0.129	11.4	79.1	0.639	0.135
21		Vase ●	11.9	67.4	0.323	0.143	16.2	68.8	0.385	0.146
22		Bottle ●	10.3	72.0	0.608	0.120	12.5	72.4	0.612	0.116
23			23.5	77.3	0.552	0.110	27.8	79.8	0.536	0.107

Table 9. The category-wise results for LASO [28] and GEAL (Ours) on the **Seen** partition of the **PIAD** dataset [63]. AUC and aIOU scores are reported in percentages (%).

#	Category	LASO [28]				GEAL (Ours)			
		aIOU↑	AUC↑	SIM↑	MAE↓	aIOU↑	AUC↑	SIM↑	MAE↓
1	Bed ●	12.0	78.0	0.469	0.126	12.8	78.4	0.473	0.120
2	Dishwasher ●	17.3	84.9	0.338	0.079	18.3	89.8	0.440	0.060
3	Laptop ●	4.5	65.4	0.192	0.122	6.3	74.5	0.201	0.100
4	Microwave ●	14.4	83.4	0.365	0.066	15.8	89.6	0.402	0.049
5	Scissors ●	3.2	66.5	0.310	0.107	3.7	69.8	0.333	0.123
6	Vase ●	5.2	58.1	0.455	0.140	6.4	54.9	0.466	0.127

Table 10. The category-wise results for LASO [28] and GEAL (Ours) on the **Unseen** partition of the **PIAD** dataset [63]. AUC and aIOU scores are reported in percentages (%).

narios, highlighting the importance of establishing robust ethical guidelines for its use.

- **Resource Intensity:** Training and deploying such sophisticated models demand significant computational resources, which can pose a challenge for smaller organizations or regions with limited access to advanced technology infrastructure. This barrier may restrict the broader adoption of the framework in resource-constrained environments.

F. Public Resource Used

In this section, we acknowledge the use of the following public resources, during the course of this work:

- LASO¹ Unknown
- IAGNet² Unknown
- PointCloud-C³ Unknown
- OOAL⁴ MIT License
- DreamGaussian⁵ MIT License
- LangSplat⁶ Gaussian-Splatting License

¹<https://github.com/y13800/LASO>

²<https://github.com/yyvhang/IAGNet>

³<https://github.com/ldkong1205/PointCloud-C>

⁴<https://github.com/Reagan1311/OOAL>

⁵<https://github.com/dreamgaussian/dreamgaussian>

⁶<https://github.com/minghanqin/LangSplat>

#	Category	LASO [28]				GEAL (Ours)				
		aIOU ↑	AUC ↑	SIM ↑	MAE ↓	aIOU ↑	AUC ↑	SIM ↑	MAE ↓	
1	StorageFurniture •	Bag •	19.8	85.4	0.535	0.085	20.6	86.7	0.572	0.084
2		Bed •	13.6	77.4	0.515	0.111	16.0	79.9	0.527	0.110
3		Bowl •	8.6	81.3	0.777	0.102	12.2	87.4	0.807	0.102
4		Clock •	6.2	84.2	0.461	0.064	9.8	84.8	0.485	0.062
5		Dishwash •	29.6	94.1	0.472	0.070	28.5	89.9	0.505	0.068
6		Display •	31.0	92.2	0.700	0.086	41.1	92.6	0.718	0.088
7		Door •	12.3	82.3	0.311	0.060	15.7	83.8	0.368	0.058
8		Earphone •	26.5	93.0	0.639	0.099	27.5	94.0	0.662	0.094
9		Faucet •	14.2	78.9	0.498	0.089	18.3	84.3	0.589	0.087
10		Hat •	3.6	67.0	0.538	0.152	9.3	72.7	0.560	0.148
11		Keyboard •	19.2	88.6	0.437	0.067	24.7	89.3	0.481	0.066
12		Knife •	12.0	89.0	0.227	0.055	12.9	87.9	0.232	0.039
13		Laptop •	14.8	91.3	0.642	0.064	22.9	93.2	0.657	0.063
14		Microwave •	28.5	95.1	0.583	0.078	29.8	95.1	0.586	0.070
15		Mug •	27.2	96.1	0.440	0.042	31.8	92.8	0.464	0.038
16		Refrigerator •	13.3	78.1	0.547	0.098	21.7	87.6	0.635	0.076
17		Chair •	25.6	92.8	0.433	0.063	24.8	93.7	0.484	0.069
18		Scissors •	28.9	89.9	0.650	0.093	28.7	89.9	0.678	0.091
19		Table •	17.5	95.4	0.661	0.053	24.9	95.9	0.684	0.045
20		TrashCan •	10.1	81.7	0.662	0.119	10.8	81.6	0.690	0.115
21		Vase •	10.9	72.1	0.323	0.137	27.8	90.4	0.499	0.100
22		Bottle •	7.9	71.1	0.630	0.125	13.5	79.5	0.650	0.116
23			20.4	81.2	0.553	0.114	28.7	81.9	0.570	0.116

Table 11. The category-wise results for LASO [28] and GEAL (Ours) on the **Seen** partition of the **LASO** dataset [28]. AUC and aIOU scores are reported in percentages (%).

#	Category	LASO [28]				GEAL (Ours)				
		aIOU ↑	AUC ↑	SIM ↑	MAE ↓	aIOU ↑	AUC ↑	SIM ↑	MAE ↓	
1	StorageFurniture •	Bag •	20.7	89.1	0.513	0.089	22.1	91.0	0.522	0.086
2		Bed •	12.2	80.6	0.553	0.115	13.6	81.4	0.563	0.113
3		Bowl •	7.5	81.3	0.744	0.125	9.1	82.5	0.749	0.119
4		Clock •	5.3	85.2	0.419	0.094	6.4	85.0	0.433	0.079
5		Dishwash •	20.7	92.4	0.443	0.069	26.0	92.4	0.470	0.065
6		Display •	23.4	86.6	0.512	0.112	25.0	87.6	0.526	0.112
7		Door •	3.4	81.3	0.324	0.095	11.7	81.4	0.355	0.066
8		Earphone •	9.5	76.8	0.454	0.130	20.8	93.5	0.639	0.091
9		Faucet •	13.8	74.1	0.442	0.098	15.1	76.8	0.470	0.095
10		Hat •	4.5	61.2	0.586	0.158	4.1	66.5	0.582	0.149
11		Keyboard •	17.9	88.1	0.422	0.069	18.3	88.3	0.423	0.067
12		Knife •	3.1	74.6	0.138	0.082	3.3	79.4	0.137	0.078
13		Laptop •	15.3	91.7	0.643	0.053	15.4	91.2	0.675	0.059
14		Microwave •	8.7	79.7	0.334	0.096	29.3	95.6	0.610	0.064
15		Mug •	11.9	90.9	0.317	0.063	14.2	91.5	0.318	0.064
16		Refrigerator •	1.7	64.5	0.381	0.174	2.5	66.6	0.511	0.157
17		Chair •	20.1	87.2	0.378	0.066	21.0	89.4	0.390	0.065
18		Scissors •	25.2	87.4	0.642	0.098	26.0	89.4	0.624	0.094
19		Table •	1.6	25.3	0.094	0.105	2.1	27.6	0.105	0.097
20		TrashCan •	7.5	70.4	0.604	0.135	7.8	72.1	0.620	0.129
21		Vase •	2.6	63.1	0.191	0.124	7.4	71.0	0.293	0.125
22		Bottle •	6.4	56.4	0.466	0.148	7.6	67.0	0.614	0.140
23			16.2	78.5	0.455	0.134	21.2	78.2	0.519	0.119

Table 12. The category-wise results for LASO [28] and GEAL (Ours) on the **Unseen** partition of the **LASO** dataset [28]. AUC and aIOU scores are reported in percentages (%).

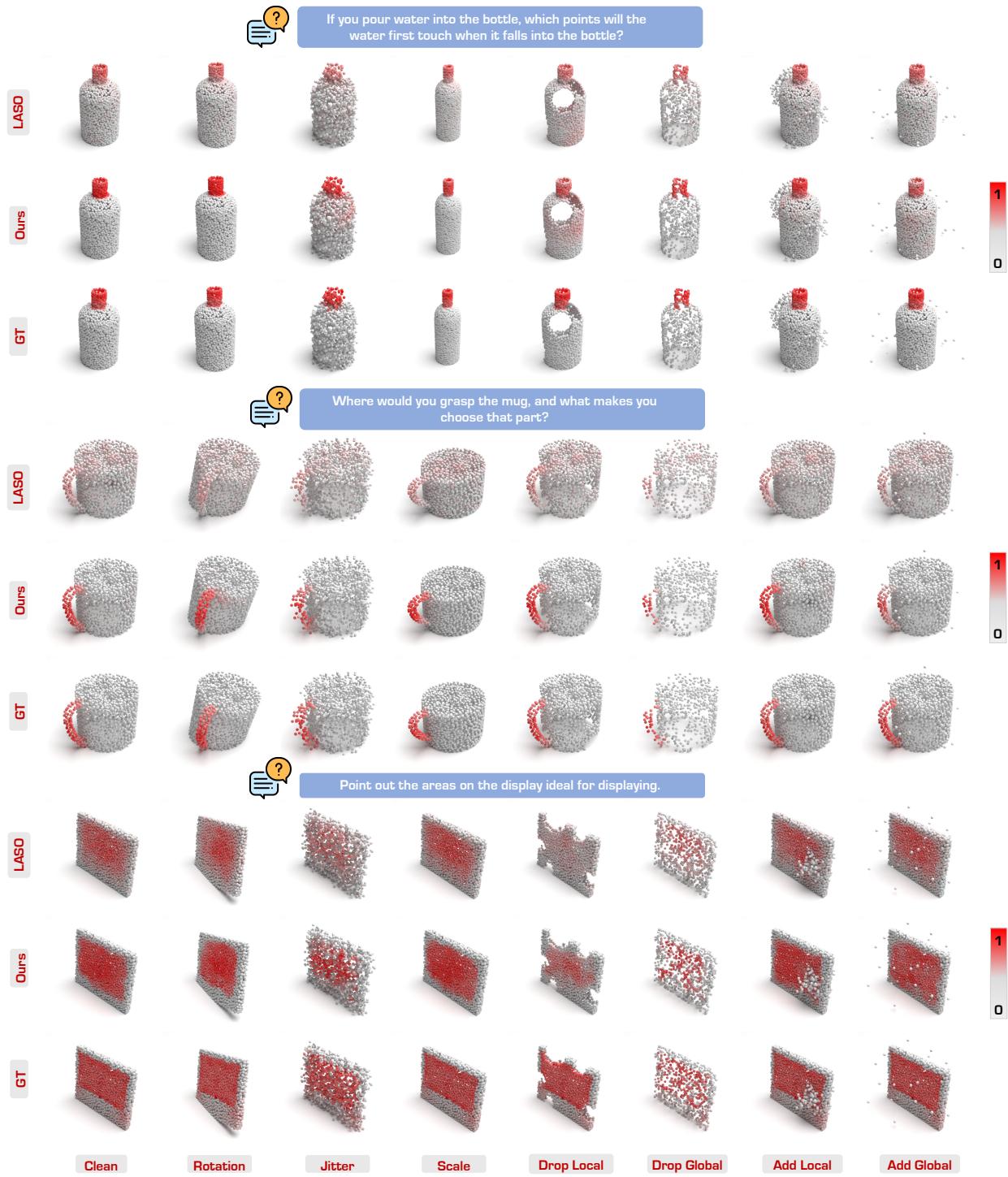


Figure 7. Qualitative comparisons between **GEAL** and **LASO** [28] on the PIAD-C dataset, highlighting the superior robustness of our method on corrupted data.

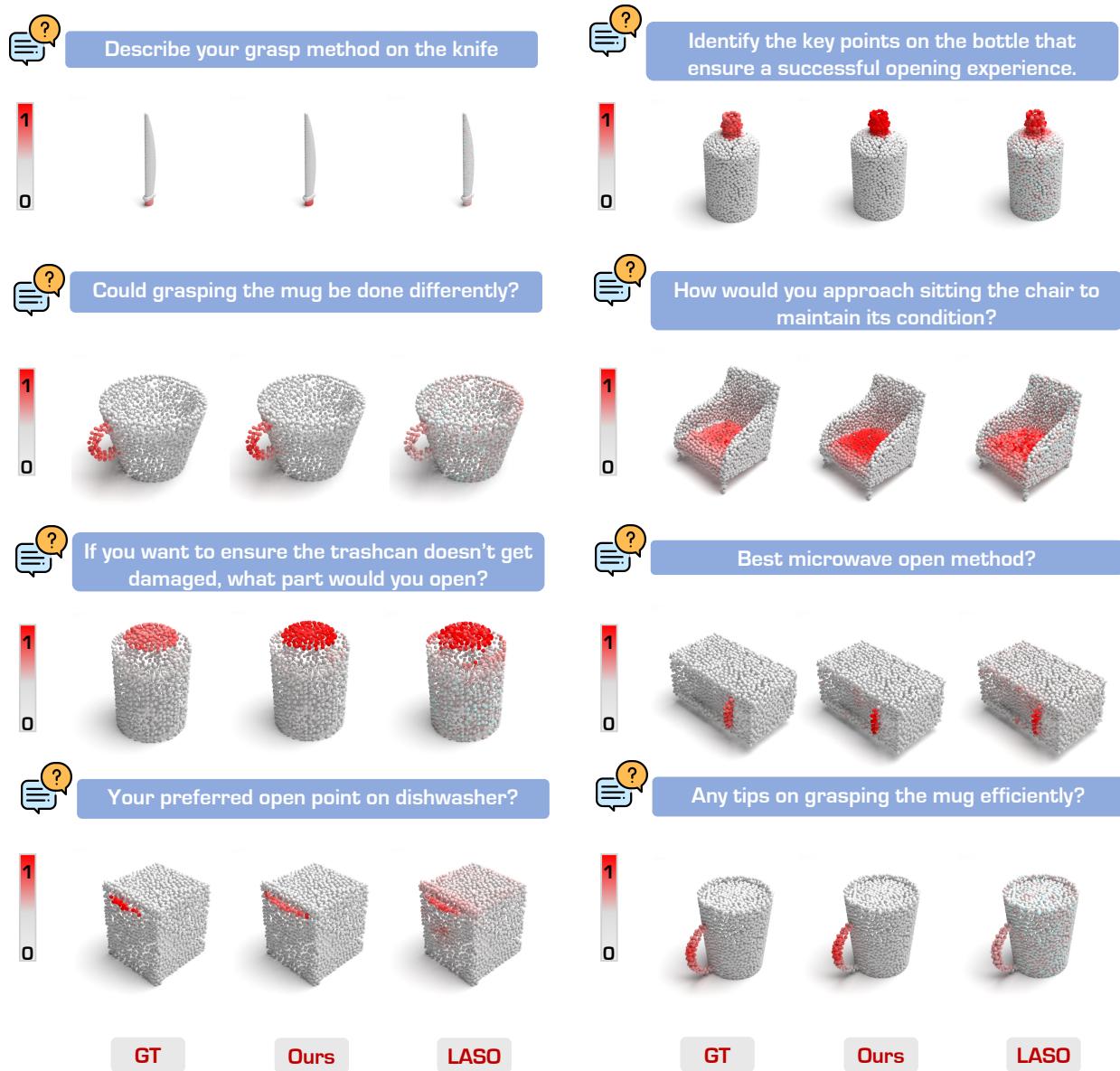


Figure 8. Qualitative comparisons between **GEAL** and **LASO** [28] on the PIAD dataset.

References

- [1] Emanuele Aiello, Diego Valsesia, and Enrico Magli. Cross-modal learning for image-guided point cloud shape completion. In *Advances in Neural Information Processing Systems*, pages 37349–37362, 2022. [6](#), [13](#), [14](#)
- [2] 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. [2](#)
- [3] Honghua Chen, Zeyong Wei, Yabin Xu, Mingqiang Wei, and Jun Wang. Imlovenet: Misaligned image-supported registration network for low-overlap point cloud pairs. In *ACM SIGGRAPH Conference Proceedings*, pages 1–9, 2022. [6](#), [13](#), [14](#)
- [4] Joya Chen, Difei Gao, Kevin Qinghong Lin, and Mike Zheng Shou. Affordance grounding from demonstration video to target image. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6799–6808, 2023. [2](#)
- [5] Ching-Yao Chuang, Jiaman Li, Antonio Torralba, and Sanja Fidler. Learning to act properly: Predicting and explaining affordances from images. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 975–983, 2018. [1](#), [2](#)
- [6] Francisco Cruz, Sven Magg, Cornelius Weber, and Stefan Wermter. Training agents with interactive reinforcement learning and contextual affordances. *IEEE Transactions on Cognitive and Developmental Systems*, 8(4):271–284, 2016. [1](#)
- [7] Leiyao Cui, Xiaoxue Chen, Hao Zhao, Guyue Zhou, and Yixin Zhu. Strap: Structured object affordance segmentation with point supervision. *arXiv preprint arXiv:2304.08492*, 2023. [1](#)
- [8] Shengcheng Deng, Xun Xu, Chaozheng Wu, Ke Chen, and Kui Jia. 3d affordancenet: A benchmark for visual object affordance understanding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1778–1787, 2021. [2](#), [11](#)
- [9] Thanh-Toan Do, Anh Nguyen, and Ian Reid. Affordancenet: An end-to-end deep learning approach for object affordance detection. In *IEEE International Conference on Robotics and Automation*, pages 5882–5889, 2018. [2](#)
- [10] Kuan Fang, Te-Lin Wu, Daniel Yang, Silvio Savarese, and Joseph J Lim. Demo2vec: Reasoning object affordances from online videos. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2139–2147, 2018. [1](#), [2](#)
- [11] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7081–7091, 2023. [2](#)
- [12] Yiran Geng, Boshi An, Haoran Geng, Yuanpei Chen, Yaodong Yang, and Hao Dong. Rlafford: End-to-end affordance learning for robotic manipulation. In *IEEE International Conference on Robotics and Automation*, pages 5880–5886, 2023. [1](#)
- [13] James J Gibson. *The ecological approach to visual perception: classic edition*. Psychology press, 2014. [1](#)
- [14] Denis Hadjivelichkov, 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](#)
- [15] Mohammed Hassanin, Salman Khan, and Murat Tahtali. Visual affordance and function understanding: A survey. *ACM Computing Surveys*, 54(3):1–35, 2021. [1](#)
- [16] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. [2](#), [9](#)
- [17] Zhi Hou, Baosheng Yu, Yu Qiao, Xiaojiang Peng, and Dacheng Tao. Affordance transfer learning for human-object interaction detection. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 495–504, 2021. [1](#)
- [18] 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. In *Conference on Robot Learning*, pages 540–562. PMLR, 2023. [2](#)
- [19] 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. [2](#), [3](#)
- [20] Sihyeon Kim, Sanghyeok Lee, Dasol Hwang, Jaewon Lee, Seong Jae Hwang, and Hyunwoo J Kim. Point cloud augmentation with weighted local transformations. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 548–557, 2021. [3](#), [9](#)
- [21] Lingdong Kong, Youquan Liu, Xin Li, Runnan Chen, Wenwei Zhang, Jiawei Ren, Liang Pan, Kai Chen, and Ziwei Liu. Robo3d: Towards robust and reliable 3d perception against corruptions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19994–20006, 2023. [2](#), [3](#), [9](#)
- [22] Hema Swetha Koppula, Rudhir Gupta, and Ashutosh Saxena. Learning human activities and object affordances from rgb-d videos. *International Journal of Robotics Research*, 32(8):951–970, 2013. [2](#)
- [23] Dogyoon Lee, Jaeha Lee, Junhyeop Lee, Hyeongmin Lee, Minhyeok Lee, Sungmin Woo, and Sangyoun Lee. Regularization strategy for point cloud via rigidly mixed sample. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15900–15909, 2021. [3](#), [9](#)
- [24] Gen Li, Varun Jampani, Deqing Sun, and Laura Sevilla-Lara. Locate: Localize and transfer object parts for weakly supervised affordance grounding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10922–10931, 2023. [1](#), [2](#), [3](#)
- [25] Gen Li, Deqing Sun, Laura Sevilla-Lara, and Varun Jampani. One-shot open affordance learning with foundation models. In *Proceedings of the IEEE/CVF Conference on Computer*

- Vision and Pattern Recognition*, pages 3086–3096, 2024. 1, 2, 3, 4
- [26] Hongxiang Li, Meng Cao, Xuxin Cheng, Yaowei Li, Zhi-hong Zhu, and Yuexian Zou. G2l: Semantically aligned and uniform video grounding via geodesic and game theory. In *IEEE/CVF International Conference on Computer Vision*, pages 12032–12042, 2023. 2
- [27] Muchen Li and Leonid Sigal. Referring transformer: A one-step approach to multi-task visual grounding. In *Advances in Neural Information Processing Systems*, pages 19652–19664, 2021. 6, 13, 14
- [28] Yicong Li, Na Zhao, Junbin Xiao, Chun Feng, Xiang Wang, and Tat-Seng Chua. Laso: Language-guided affordance segmentation on 3d object. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14251–14260, 2024. 1, 2, 5, 6, 7, 8, 10, 11, 13, 14, 15, 16, 17, 18
- [29] Chang Liu, Henghui Ding, and Xudong Jiang. Gres: Generalized referring expression segmentation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23592–23601, 2023. 6, 13, 14
- [30] Liu Liu, Wenqiang Xu, Haoyuan Fu, Sucheng Qian, Qiaojun Yu, Yang Han, and Cewu Lu. Akb-48: A real-world articulated object knowledge base. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14809–14818, 2022. 2
- [31] Shaowei Liu, Subarna Tripathi, Somdeb Majumdar, and Xiaolong Wang. Joint hand motion and interaction hotspots prediction from egocentric videos. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3282–3292, 2022. 2
- [32] Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 364, 2019. 6
- [33] Jorge M Lobo, Alberto Jiménez-Valverde, and Raimundo Real. Auc: a misleading measure of the performance of predictive distribution models. *Global Ecology and Biogeography*, 17(2):145–151, 2008. 7, 13
- [34] Liangsheng Lu, Wei Zhai, Hongchen Luo, Yu Kang, and Yang Cao. Phrase-based affordance detection via cyclic bilateral interaction. *IEEE Transactions on Artificial Intelligence*, 4(5):1186–1198, 2023. 2
- [35] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. One-shot affordance detection. *arXiv preprint arXiv:2106.14747*, 2021. 2
- [36] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. Grounded affordance from exocentric view. *arXiv preprint arXiv:2208.13196*, 2022. 2
- [37] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. Learning affordance grounding from exocentric images. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2252–2261, 2022. 1, 2
- [38] Junyu Luo, Jiahui Fu, Xianghao Kong, Chen Gao, Haibing Ren, Hao Shen, Huaxia Xia, and Si Liu. 3d-sps: Single-stage 3d visual grounding via referred point progressive selection. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16454–16463, 2022. 6, 13, 14
- [39] Jinpeng Mi, Song Tang, Zhen Deng, Michael Goerner, and Jianwei Zhang. Object affordance based multimodal fusion for natural human-robot interaction. *Cognitive Systems Research*, 54:128–137, 2019. 2
- [40] Jinpeng Mi, Hongzhuo Liang, Nikolaos Katsakis, Song Tang, Qingdu Li, Changshui Zhang, and Jianwei Zhang. Intention-related natural language grounding via object affordance detection and intention semantic extraction. *Frontiers in Neurorobotics*, 14:26, 2020. 2
- [41] Kaichun Mo, Shilin Zhu, Angel X. Chang, Li Yi, Subarna Tripathi, Leonidas J. Guibas, and Hao Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 909–918, 2019. 2
- [42] Kaichun Mo, Yuzhe Qin, Fanbo Xiang, Hao Su, and Leonidas Guibas. O2o-afford: Annotation-free large-scale object-object affordance learning. In *Conference on Robot Learning*, pages 1666–1677. PMLR, 2022. 2
- [43] Austin Myers, Ching L Teo, Cornelia Fermüller, and Yiannis Aloimonos. Affordance detection of tool parts from geometric features. In *IEEE International Conference on Robotics and Automation*, pages 1374–1381, 2015. 2
- [44] Tushar Nagarajan, Christoph Feichtenhofer, and Kristen Grauman. Grounded human-object interaction hotspots from video. In *IEEE/CVF International Conference on Computer Vision*, pages 8688–8697, 2019. 2
- [45] Toan Nguyen, Minh Nhat Vu, An Vuong, Dzung Nguyen, Thieu Vo, Ngan Le, and Anh Nguyen. Open-vocabulary affordance detection in 3d point clouds. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5692–5698, 2023. 2
- [46] Maxime Oquab, Timothée Darct, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2, 3, 4, 5, 6
- [47] 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. 4, 5
- [48] 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. 2, 3
- [49] Md Atiqur Rahman and Yang Wang. Optimizing intersection-over-union in deep neural networks for image segmentation. In *International Symposium on Visual Computing*, pages 234–244, 2016. 7, 13
- [50] Jiawei Ren, Lingdong Kong, Liang Pan, and Ziwei Liu. Benchmarking and analyzing point cloud robustness under corruptions. *Preprint*, 2022. 2, 6
- [51] Jiawei Ren, Liang Pan, and Ziwei Liu. Benchmarking and analyzing point cloud classification under corruptions. In *International Conference on Machine Learning*, pages 18559–18575. PMLR, 2022. 2, 9

- [52] Anirban Roy and Sinisa Todorovic. A multi-scale cnn for affordance segmentation in rgb images. In *European conference on computer vision*, pages 186–201. Springer, 2016. 2
- [53] Michael J Swain and Dana H Ballard. Color indexing. *International Journal of Computer Vision*, 7(1):11–32, 1991. 7, 13
- [54] Xun Tan, Xingyu Chen, Guowei Zhang, Jishiyu Ding, and Xuguang Lan. Mbdf-net: Multi-branch deep fusion network for 3d object detection. In *International Workshop on Multimedia Computing for Urban Data*, pages 9–17, 2021. 6, 13, 14
- [55] Spyridon Thermos, Petros Daras, and Gerasimos Potamianos. A deep learning approach to object affordance segmentation. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2358–2362. IEEE, 2020. 2
- [56] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34:24261–24272, 2021. 5
- [57] Jie Wang, Lihe Ding, Tingfa Xu, Shaocong Dong, Xinli Xu, Long Bai, and Jianan Li. Sample-adaptive augmentation for point cloud recognition against real-world corruptions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14330–14339, 2023. 3, 9
- [58] Jie Wang, Tingfa Xu, Lihe Ding, and Jianan Li. Target-guided adversarial point cloud transformer towards recognition against real-world corruptions. *arXiv preprint arXiv:2411.00462*, 2024. 3, 9
- [59] Cort J Willmott and Kenji Matsuura. Advantages of the mean absolute error (mae) over the root mean square error (rmse) in assessing average model performance. *Climate Research*, 30(1):79–82, 2005. 7, 13
- [60] Chao Xu, Yixin Chen, He Wang, Song-Chun Zhu, Yixin Zhu, and Siyuan Huang. Partafford: Part-level affordance discovery from 3d objects. *arXiv preprint arXiv:2202.13519*, 2022. 2
- [61] Danfei Xu, Dragomir Anguelov, and Ashesh Jain. Pointfusion: Deep sensor fusion for 3d bounding box estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 244–253, 2018. 6, 13, 14
- [62] Xinli Xu, Shaocong Dong, Lihe Ding, Jie Wang, Tingfa Xu, and Jianan Li. Fusionrcnn: Lidar-camera fusion for two-stage 3d object detection. *Remote Sensing*, 15(7):1839, 2023. 6, 13, 14
- [63] Yuhang Yang, Wei Zhai, Hongchen Luo, Yang Cao, Jiebo Luo, and Zheng-Jun Zha. Grounding 3d object affordance from 2d interactions in images. In *IEEE/CVF International Conference on Computer Vision*, pages 10905–10915, 2023. 1, 2, 6, 7, 8, 10, 11, 12, 13, 14, 15
- [64] Xue Zhao, Yang Cao, and Yu Kang. Object affordance detection with relationship-aware network. *Neural Computing and Applications*, 32(18):14321–14333, 2020. 2
- [65] Zihan Zhong, Zhiqiang Tang, Tong He, Haoyang Fang, and Chun Yuan. Convolution meets lora: Parameter efficient finetuning for segment anything model. *arXiv preprint arXiv:2401.17868*, 2024. 4
- [66] Zhuangwei Zhuang, Rong Li, Kui Jia, Qicheng Wang, Yuanqing Li, and Mingkui Tan. Perception-aware multi-sensor fusion for 3d lidar semantic segmentation. In *IEEE/CVF International Conference on Computer Vision*, pages 16280–16290, 2021. 6, 13, 14