# WorldAfford: Affordance Grounding based on Natural Language Instructions

Changmao Chen[1], Yuren Cong[2], Zhen Kan[1]

[1] University of Science and Technology of China
[2] TNT, Leibniz University Hannover
{abc,lncs}@uni-heidelberg.de, cong@tnt.uni-hannover.de

**Abstract.** Affordance grounding aims to localize the interaction regions for the manipulated objects in the scene image according to given instructions, which is essential for Embodied AI and manipulation tasks. A critical challenge in affordance grounding is that the embodied agent should understand human instructions and analyze which tools in the environment can be used, as well as how to use these tools to accomplish the instructions. Most recent works primarily supports simple action labels as input instructions for localizing affordance regions, failing to capture complex human objectives. Moreover, these approaches typically identify affordance regions of only a single object in object-centric images, ignoring the object context and struggling to localize affordance regions of multiple objects in complex scenes for practical applications. To address this concern, for the first time, we introduce a new task of affordance grounding based on natural language instructions, extending it from previously using simple labels for complex human instructions. For this new task, we propose a new framework, **WorldAfford**. We design a novel Affordance Reasoning Chain-of-Thought Prompting to reason about affordance knowledge from LLMs more precisely and logically. Subsequently, we use SAM and CLIP to localize the objects related to the affordance knowledge in the image. We identify the affordance regions of the objects through an affordance region localization module. To benchmark this new task and validate our framework, an affordance grounding dataset, LLMaFF, is constructed. We conduct extensive experiments to verify that WorldAfford performs state-of-the-art on both the previous AGD20K and the new LLMaFF dataset. In particular, WorldAfford can localize the affordance regions of multiple objects and provide an alternative when objects in the environment cannot fully match the given instruction. The code will be released after the publication of this work.

**Keywords:** Affordance Grounding · Natural Language Instruction · LLM

## 1 Introduction

Embodied agents can interact with a physical environment and potentially perform heavy tasks based on human instructions. In order for robots to better

manipulate objects in complex scenes, it is urgent to understand which part of the object is the interaction region. Affordance grounding, which aims to localize potential interaction regions for the manipulated objects in the scene image depending on the given instruction, can provide a new experience for Embodied AI and has the potential to significantly increase efficiency and flexibility. As a result, it has recently attracted a significant amount of attention [15,22,32,33,49].
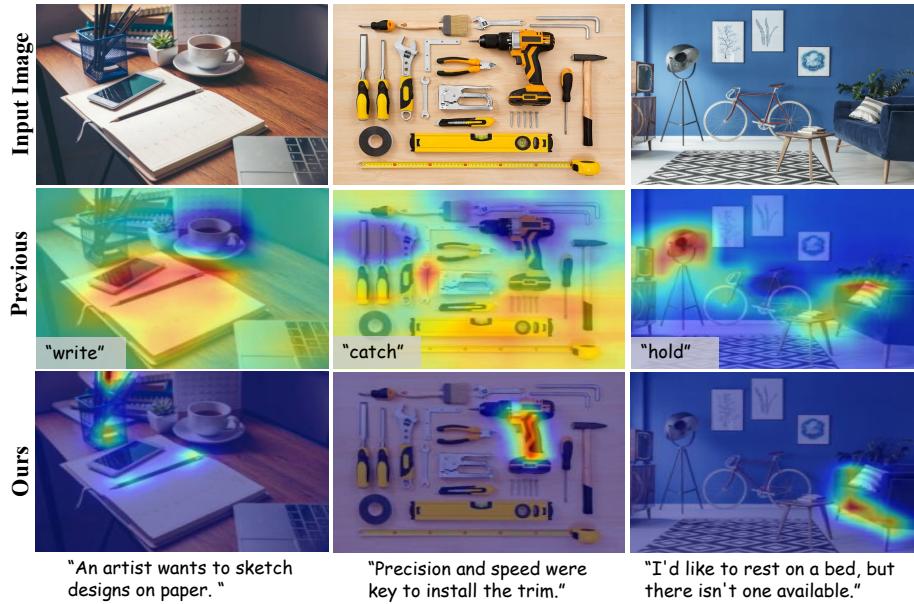


**Fig. 1:** Different from previous works only using naïve action labels for affordance grounding, WorldAfford can derive affordance knowledge from LLMs and precisely localize the affordance regions corresponding to natural language instructions. In this way, our framework can work effectively in complex open-world environments. The results in the second row are from Cross-view-AG+ [33].

A critical challenge in affordance grounding is instruction comprehension, which means that the embodied agent should understand the human instructions and reason about the actions it is going to perform, which emphasizs active interaction between humans and their environment rather than passive detection. Furthermore, the agent should analyze which tools in the usage environment can accomplish the given instructions and localize the interaction regions (*i.e.,* affordance regions) on the objects. These challenges are expected to be alleviated through using large-scale vision-language foundation models. Unfortunately, the currently available models [8, 24, 27, 42, 50, 54] have not performed satisfactorily on this particular task.

Most recent works [22,32,33,49] attempt to transfer knowledge from exocentric images of an object in an active state to egocentric images where the ob-

ject is not being used. They have achieved impressive progress, making dataset collection easier and learning that the affordance region of an object changes dynamically depending on the different given instructions. Nevertheless, current approaches can only support simple action labels (*e.g.,* "catch" shown in Fig. 1) as input instructions, which cannot express complex human goals. Besides, these methods can only identify the affordance region of a single object in object-centric images, overlook object context, and still fall short in localizing the affordance regions of multiple objects in complex scene images for practical applications in the real world. In this paper, for the first time, we introduce a new task of affordance grounding based on natural language instructions, extending affordance grounding from previously using simple action labels to complex natural language instructions. This new task moves toward real-world applications with significant implications on Embodied AI. For this task, we propose a novel framework, WorldAfford, which integrates the large language model (LLM), Segment Anything model (SAM) [21], and CLIP [42]. We first use the LLM to process the natural language instruction. To reason about affordance knowledge from the LLM more precisely and logically, we design a novel Affordance Reasoning Chain-of-Thought Prompting (ARCoT) including Object-Oriented Reasoning Prompting and Action-Oriented Reasoning Prompting. Subsequently, we employ SAM and CLIP to segment and select the objects associated with the actions inferred by the LLM. Moreover, a Weighted Context Broadcasting module (WCB) is proposed and integrated into the affordance region localization module. It allows our framework to focus on more informative objects and to identify affordance regions of multiple objects. To benchmark the new task and validate our framework, we constructed a new dataset, LLMaFF, containing 550 test images with natural language instructions and manually labeled affordance maps. Experimental results demonstrate that our framework outperforms the previous methods both on the existing AGD20K [32] dataset and the new LLMaFF dataset. Our main contributions can be summarized as follows:

- We introduce a new task of affordance grounding based on natural language instructions, extending affordance grounding from using simple action labels to complex natural language instructions.

- We propose a framework for this new task named WorldAfford, which integrates the LLM and other vision models. To reason about affordance knowledge from LLMs, we introduce an Affordance Reasoning Chain-of-Thought Prompting. In addition, we propose a Weighted Context Broadcasting module, allowing WorldAfford to localize affordance regions of multiple objects.

- A new dataset LLMaFF is constructed to benchmark the new task.

- We conduct extensive experiments to validate that our model performs state-of-the-art on both the AGD20K dataset and our new LLMaFF dataset.

## 2   Related Work

**Affordance Grounding.** Visual affordance grounding has been intensively explored in the fields of robotics and computer vision [9, 15, 25, 31, 34–36, 46, 47, 51–53]. Traditional approaches [6, 10, 11, 29, 39] mainly learn the affordance through fully supervised learning. Nagaragan *et al.* [37] learn affordance knowledge by watching human-object interaction videos. Luo *et al.* [32] propose a Cross-view-AG knowledge transfer framework for affordance grounding, in which the affordance knowledge is acquired from exocentric human-object interactions, and transfer to egocentric images. Li *et al.* [22] extract object-related information from exocentric images and match it to the objects to localize the affordance regions. However, such methods use only naive action labels for affordance grounding, which may not meet the requirements of practical applications. In this work, we use flexible natural language as supervision to guide agents in localizing affordance regions of multiple objects in complex scenes images.

**Large language models and Chain-of-Thought.** Large language models (LLMs), with their extensive world knowledge, play a central role in enabling embodied agents to interpret and execute tasks from human natural language instructions. While some studies [2, 7, 17, 18]have primarily used LLMs to guide object grasping with robotic arms, focusing on basic object perception without considering the fine-grained shapes, functions, or uses of the objects. Our work differs by using LLMs and the affordance reasoning Chain-of-Thought (ARCoT) method to interpret open-world human instructions and reason about a wide range of objects in the environment. Recent studies [12, 26, 48, 50] find that CoT can dramatically improve the performance of LLMs, particularly when dealing with complex tasks involving reasoning, demonstrate that CoT significantly improves the performance of LLMs in complex reasoning tasks. To the best of our knowledge, we are the first to explore CoT for affordance grounding.

**Vision Foundation Model for Affordance Grouding** Vision language models have shown promising results in robotics applications [16, 20]. Some works [40, 45] use vision-language models for affordance detection in 3D point clouds. They focus on the different affordances of individual objects, do not include human instructions, cannot generalize to unseen objects, and require extensive manual annotation. Li *et al.* [23] propose a vision-language framework to address the problem of one-shot affordance learning. Ren *et al.* [43] combines GroundingDINO [28] and SAM to segment objects in image based on object names. Luddecke *et al.* [30] design CLIPSeg using the CLIP as a backbone, expanded with a Transformer-based decoder, to enable segmentation based on image or text inputs. In contrast to the above work, we use CLIP for semantic understanding and SAM for spatial understanding to select objects related to language instruction. This combination allows for the accurate identification of objects as dictated by textual input.

**Affordance Grounding Dataset** Affordance grounding [5, 11, 22, 32–34, 37] has traditionally focused on datasets such as AGD20K [32] and OPRA [11] mainly for

single actionable object scenarios. Recently, Hadjivelichkov *et al.* [14] introduce the UMD-i dataset, which focuses on single objects and uses pixel-level labels for training, mainly for one-shot affordance learning. Nguyen *et al.* [38] propose the IIT-AFF dataset, which assigns affordance labels to each pixel in its images, lacks semantic information about affordance and uses only image inputs for supervised learning. To address these limitations, we present the LLMaFF dataset, which contains scenes with multiple objects and complex language instructions.

## 3    Task Definition and LLMaFF Dataset

Given an image $I$ and a natural language instruction $t$, affordance grounding based on natural language instructions aims to localize the interaction regions of objects in the scene image and the instruction can be completed through these interactions. Compared to the setting in previous works [14, 22, 32, 33], affordance grounding based on natural language instructions is more oriented towards practical applications in the real world since there is no restriction on the number of objects in the image and complexity of the input instructions.

   To facilitate and benchmark this new task, we construct a new dataset, LL-MaFF, consisting of 550 complex environmental images with natural language instructions and manually labelled affordance maps. The data collection pipeline is shown in Fig. 2. Annotators need to collect real-world images, create instructions that express the purpose of interacting with objects in the environment, and annotate ground truth (GT) based on these instructions. The source images of our dataset are primarily sampled from IIT-AFF [38]. Due to the limited object categories of IIT-AFF, we augment the dataset with the images sampled from Ego4D [13] and the Internet.  AGD20K [32] annotates the affordance re-
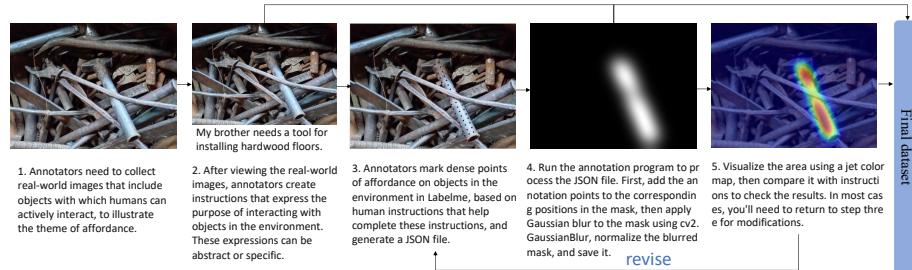


**Fig. 2:** Data Collection Pipeline for our WorldAfford benchmark.

gions with sparse points and applies a Gaussian kernel to generate ground truth. In contrast, we employ dense points to annotate the affordance map of multiple objects based on the language instructions, which requires careful identification of the objects and their interactions. We find that the density and distribution of
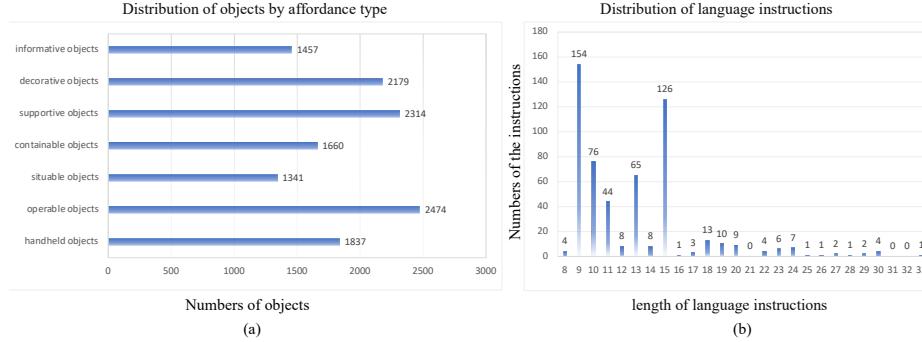
**Fig. 3:** Statistics of objects and instructions in WorldAfford. (a) Distribution of objects by affordance type. (b) Distribution of language instructions

the points have a significant impact on the labelling results, thus ensuring a uniform distribution of annotation points across multiple objects is crucial to avoid certain regions in the affordance map appearing blank or with faint heat. Based on the affordances of objects in the environment, we categorized them into eight types: handheld objects(1837), operable objects(2474), situable objects(1341), containable objects(1660), supportive objects(2314), decorative objects(2179), and informative objects(1457). We also conducted a statistical analysis of the length of human language instructions in the dataset, as shown in Fig. 3. We demonstrate some examples from the LLMaFF dataset in Fig. 4.
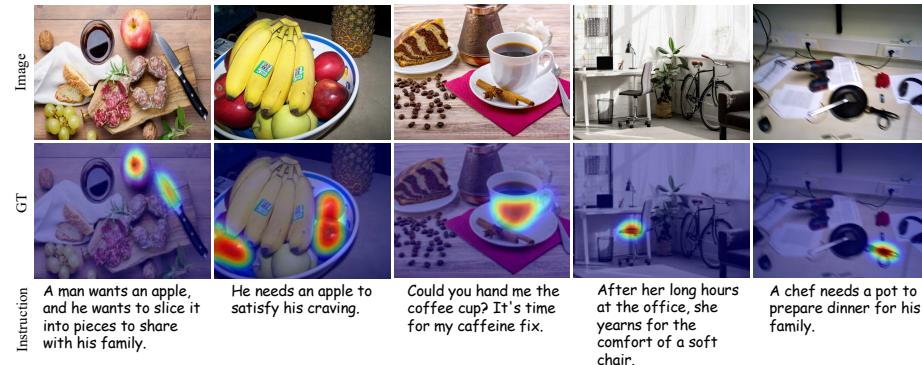


**Fig. 4:** Examples from the LLMaFF dataset. The first row presents the images from diverse environments. The second row shows the ground-truth affordance maps, and the third row shows the corresponding natural language instructions.
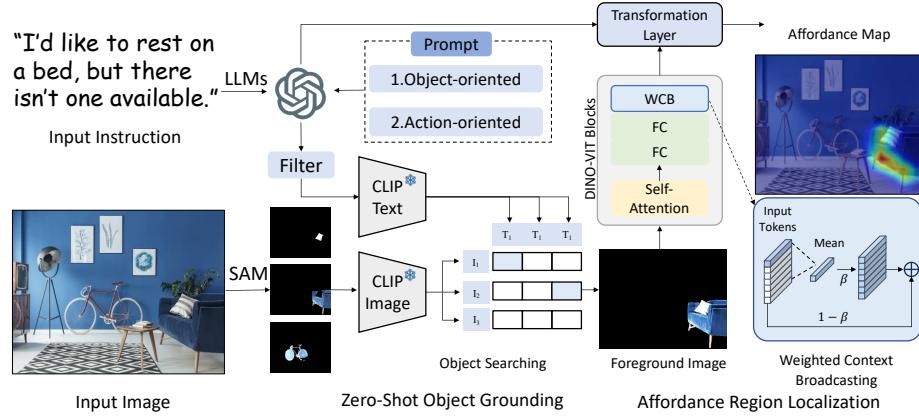
**Fig. 5:** Framework of WorldAfford. We first derive affordance knowledge from LLMs with the proposed affordance reasoning chain-of-thought. Subsequently, we utilize SAM and CLIP for zero-shot object grounding. Moreover, we design a WCB module and integrate it into the affordance region localization module to localize the affordance regions of multiple objects.

## 4   WorldAfford Framework

We propose WorldAfford as a general framework for affordance grounding incorporating complex instruction understanding and multi-object affordance localization at a very low training cost. We first use the LLM [1] to analyze the given instruction and derive affordance knowledge through the affordance reasoning chain-of-thought prompting. Subsequently, we utilize SAM [21] and CLIP [42] to implement zero-shot multi-object grounding, segmenting and selecting the objects associated with the sub-actions provided by the LLM. Furthermore, we design a Weighted Context Broadcasting (WCB) and integrate it into the affordance region localization module to localize the affordance regions of multiple objects. The overall framework of WorldAfford is illustrated in Fig. 5.

### 4.1   Affordance Reasoning Chain-of-Thought Prompting

When humans receive an instruction, they initially consider which tools (objects) might facilitate the task and decompose the complex instruction into a series of simpler actions for execution. In this paper, we prompt the LLM to mimic this aspect of human cognitive processing, and reason about relevant affordance information using the extensive world knowledge learned from large-scale data.

Rather than relying on direct inference, our approach employs a straightforward and effective chain-of-thought prompting to enhance the capabilities of LLMs in affordance reasoning. The proposed chain-of-thought prompting consists of two primary strategies: (1) **object-oriented reasoning prompting**,

and (2) **action-oriented reasoning prompting**. It enables the LLM to derive crucial object and action information from natural language instructions. Furthermore, our method allows the LLM to identify alternative tools when the optimal one is not available, demonstrating its adaptability to various scenarios.

**Object-Oriented Reasoning Prompting.** We first utilize the LLM to reason about the possible objects that can afford the given instruction. Considering that multiple objects are likely to be necessary for completing the instruction, the LLM is requested to output a set of object categories $\mathcal{O}$:

$$\mathcal{O} = LLM(k, t, p_{obj}), \tag{1}$$

where $k$ indicates the size of the object set and $t$ denotes the given natural language instruction. The prompt $p_{obj}$ for object-oriented reasoning is specifically designed as follows:

*Prompt: What are the [#k] most common objects that can be used if [#t]?*
*Output: Chair…, Hammock…, Blanket and Pillows…*

The object-oriented reasoning prompts the LLM to provide diverse objects suitable for an action. Moreover, it associates alternative tools in case the best tool does not exist in the environment, which facilitates the accomplishment of the instruction. These inferred object categories from the LLM are further utilized for subsequent action-oriented reasoning.

We designed a filter function based on the large model's output about object descriptions to filter out excessive explanatory text, which sometimes includes irrelevant objects, hindering the subsequent object search.

**Action-Oriented Reasoning Prompting.** Different from previous work on locating the affordance for single action labels, to address the complex natural language instructions, we utilize the powerful prior knowledge of the LLM to decompose a complex instruction into several simple sub-actions. Given a pre-defined predicate list, we prompt the LLM to select the appropriate predicates from it and assign these predicates to the objects in the object set $\mathcal{O}$. A set of sub-actions $\mathcal{A}$ is generated as follows:



**Fig. 6:** The affordance reasoning Chain-of-Thought prompting.

$$\mathcal{A} = \mathrm{LLM}(\mathcal{O}, l_p, t, p_{act}), \tag{2}$$

where $l_p$ indicates the pre-defined predicate list and $t$ denotes the given instruction. Each sub-action consists of a predicate and an object. The prompt $p_{act}$ for action-oriented reasoning is specifically designed as follows:

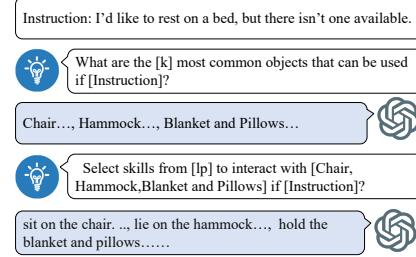*Prompt: Select skills from [#$l_p$] to interact with the above [#$\mathcal{O}$] to help me if [#t]?*

*Output: sit on the chair..., lie on the hammock..., hold the blanket and pillows...*

The inferred sub-actions from the LLM are further utilized as the input of the following affordance region localization module.

We use the LLM to extract object-level knowledge and aggregate action-level knowledge. In the inference process of the LLM, irrelevant information in the natural language instruction is ignored and the highly abstract instruction is transformed into a series of executable sub-actions. The powerful reasoning ability and adaptive results of the affordance reasoning chain-of-thought facilitate the subsequent zero-shot object grounding and the affordance region localization.

### 4.2   Zero-shot Multiple Object Grounding

To integrate the affordance knowledge provided by the LLM with visual information about the environment, our approach leverages the capabilities of Segment Anything Model (SAM) [21] and CLIP [42] to effectively ground the relevant objects in the scene image according to the given natural language instruction. The impressive zero-shot performance of SAM and CLIP enables our framework to precisely localize objects across the open world without the need for extensive and expensive training on large-scale datasets.

Initially, SAM produces $N$ segmentation masks for the input image. These masks, while precisely segmented, lack semantic labels and unavoidably contain irrelevant objects. In order to obtain the object masks that are relevant to the given instruction, CLIP is integrated to compute the similarity between the visual appearance of the masks and the object categories provided by the LLM. We extract the corresponding regions from the original image $\boldsymbol{I}$ based on the segmentation masks. Subsequently, the cropped regions $\boldsymbol{m}$ are encoded by the CLIP image encoder $\mathrm{E}_{image}$ while the textual object categories $\boldsymbol{o}$ are encoded by the CLIP text encoder $\mathrm{E}_{text}$. The probability $p$ of the mask being classified as the $i$-th object category can be formulated as:

$$p = \frac{\exp(\mathrm{sim}(\mathrm{E}_{image}(\boldsymbol{m}), \mathrm{E}_{text}(\boldsymbol{o}_i))/\alpha)}{\sum_{o_i \in \mathcal{O}} \exp(\mathrm{sim}(\mathrm{E}_{image}(\boldsymbol{m}), \mathrm{E}_{text}(\boldsymbol{o}_i))/\alpha)}, \tag{3}$$

where $\mathrm{sim}(,)$ denotes the cosine similarity function and $\mathcal{O}$ indicates the set of object categories from the LLM. The scaling factor $\alpha$ is set to 0.1 in practice. We establish a boundary to determine whether the masks from SAM are valid. The masks with probability $p$ above the boundary are identified as valid masks. With these active masks, we construct a full-view segmentation mask in which the region covered by the valid masks is viewed as foreground, while the remaining area is considered as background. This full-view mask is the same size as the input image and is further used for affordance region localization.

### 4.3   Affordance Region Localization

To localize the affordance region of the objects in the image corresponding to the given instruction, we employ LOCATE [22] and enhance the grounding performance through two crucial improvements: (1) We use the full-view mask resulting from zero-shot multi-object grounding to preserve the foreground and mask

off the background as the input, rather than the entire image. (2) We propose a weighted context broadcasting (WCB) module, seamlessly integrating it into DINO-ViT [4] to enable the model to prioritize informative objects. With these improvements, our approach outperforms the original LOCATE and can localize multiple affordance regions with the knowledge provided by the LLM.

We utilize the full-view mask from zero-shot multi-object grounding to mask off the irrelevant objects in the image. The relevant objects are preserved and the image is forwarded into DINO-ViT to extract deep part-aware features. Different from LOCATE, we design a Weighted Context Broadcasting (WCB) module inspired by Context Broadcasting [19] and incorporate it into DINO-ViT as demonstrated in Fig. 5. Given a sequence of $N$ patch tokens, the WCB module combines the average context tokens with the input tokens in a weighted manner as follows:

$$\text{WCB}(x_i) = x_i * \beta + \frac{1}{N} \sum_{j=1}^{N} x_j * (1 - \beta), \tag{4}$$

where the weight $\beta$ is an empirically determined hyperparameter. In order to improve the model's capability to perceive multiple objects, it is expected that the attention maps in the self-attention modules of DINO-ViT are dense rather than sparse. It has been discussed in [19] that aggregating the average context token can facilitate the self-attention modules to learn dense attention maps. However, such simple aggregation makes training difficult since the target attention is unknown and uncertain. To solve this issue, we introduce a weight to balance the aggregation. With the proposed WCB, the target attention is easier to learn and the model can focus on more informative objects. The experiment in Sec. 5.3 also demonstrates that our approach outperforms the previous works [33] in terms of affordance grounding for objects.

The feature maps generated by DINO-ViT are further refined by a transformation layer including a feed-forward layer and two subsequent convolutional layers. We follow the training strategy of LOCATE [22] to transfer affordance knowledge from exocentric images to egocentric images. To predict the affordance maps, a convolutional layer with a window size of $1 \times 1$ is utilized to project the channel number to the total number of the action categories in the pre-defined predicate list $l_p$. We aggregate the affordance maps corresponding to the action categories provided by the LLM and normalized them to limit the activation values in the map between 0 and 1 as the final output.

## 5    Experiments

### 5.1    Datasets and Evaluation Metrics

We conduct experiments on the AGD20K [32] dataset and our proposed LLMaFF dataset. AGD20K dataset stands out as the only large-scale dataset containing 20,061 demonstration images and 6,060 object images for training. We select the test set of 1,675 images to assess the performance of our method in the

**Table 1:** Comparison of WorldAfford with other state-of-the-art affordance grounding methods on AGD20K. The best numbers are highlighted in **bold**.

| Approach | Input Instruction | KLD↓ | SIM↑ | NSS↑ |
|---|---|---|---|---|
| Hotspots [37] | Action Label | 1.773 | 0.278 | 0.615 |
| Cross-view-AG [32] | Action Label | 1.538 | 0.334 | 0.927 |
| Affcorr [14] | Action Label | 1.407 | 0.359 | 1.026 |
| LOCATE [22] | Action Label | 1.226 | 0.401 | 1.177 |
| Cross-view-AG+ [33] | Action Label | 1.213 | 0.403 | 1.242 |
| WorldAfford(ours) | Action Label | **1.201** | **0.406** | **1.255** |

**Table 2:** Comparison on LLMaFF dataset. We manually select labels for the other methods to comparison with them. WorldAfford outperforms all previous methods across all evaluation metrics. The best results are highlighted in **bold**.

| Approach | Input Instruction | KLD↓ | SIM↑ | NSS↑ |
|---|---|---|---|---|
| Cross-view-AG+ [33] | Action Label | 2.927 | 0.123 | -0.194 |
| Cross-view-AG [32] | Action Label | 2.887 | 0.119 | 0.118 |
| LOCATE [22] | Action Label | 1.958 | 0.212 | 1.713 |
| WorldAfford(ours) | Natural Language | **1.163** | **0.386** | **2.819** |

task of affordance grounding guided by a single action label. LLMaFF dataset includes 550 complex environment images with natural language instructions and affordance maps. We conduct experiments on the LLMaFF dataset to evaluate the performance of our method in the task of affordance grounding based on natural language instructions. We use KLD [3], SIM [44], and NSS [41] as metrics to measure the correspondence between the predicted affordance map and the ground truth. Specifically, KLD quantifies the divergence between the predicted affordance maps and the distribution of ground truth images, SIM assesses the similarity between the affordance maps and the ground truth, and NSS measures the agreement between the prediction and the ground truth.

**Training strategy and comparison fairness** Only the affordance region localization module requires training, while other modules are frozen. The training is only performed on AGD20K with the same settings as the baselines at a low trainging cost. Overall, we keep the comparison as fair as possible.

### 5.2  Implementation Details.

We use GPT-4 [1] as the large language model, while both the CLIP [42] and the Segment Anything Model (SAM) [21] implement object matching and segmentation in a zero-shot fashion. The affordance information is extracted from the output of the large language model by removing most of the irrelevant text to allow the CLIP to more accurately localize the position of objects. The affordance region localization module is trained on a RTX 3090 GPU. We load the pre-trained DINO-ViT [4] model and finetune the features it extracts from
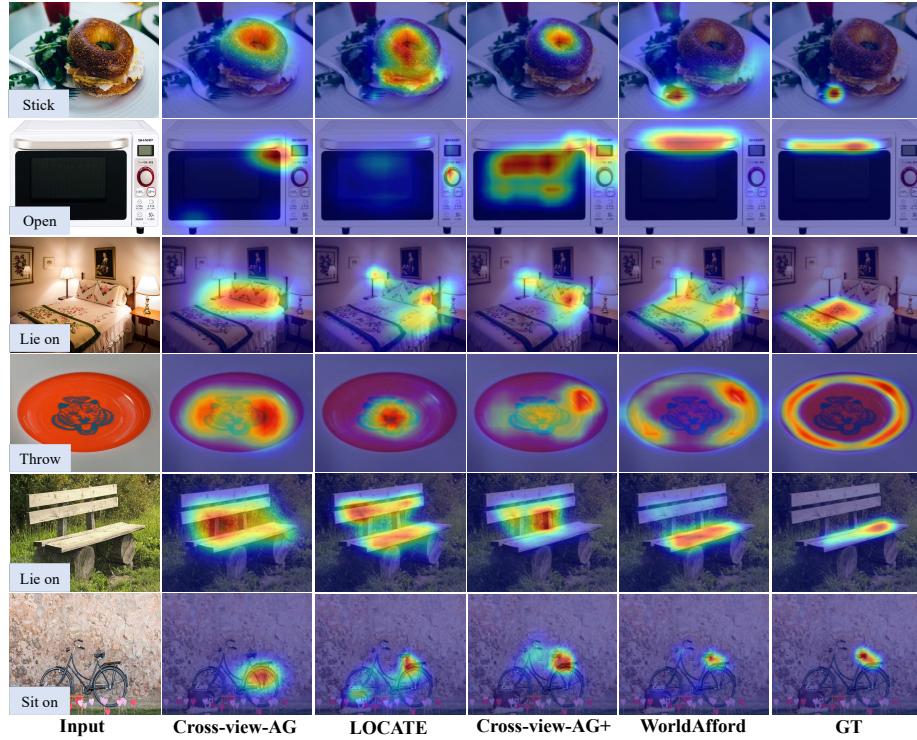
**Fig. 7:** Visual comparison on the AGD20K dataset. Compared to previous methods, our method can infer more precise affordance maps.

images. We set the weight $\beta$ to 0.88, and the number $k$ in Eq. (1) is set to 3. We use a learning rate of 0.005, a decay factor of 5e-4, a batch size of 16, and train the affordance region localization module for 35 epochs.

## 5.3 Quantitative results

For comparison with the previous affordance grounding approaches, we first validate our framework on the AGD20K dataset. AGD20K is widely used in the affordance grounding approaches [22, 32, 33], which employ a simple action label to localize the affordance regions of a single object in the object-centric images. Since WorldAfford uses natural language instructions as the input, a direct comparison is difficult to achieve. To solve this problem, we only use action labels as input to the affordance region localization module and compare it to these approaches. The results in Tab. 1 show that even in this simplified setting, our approach still outperforms the previous methods. WorldAfford establishes a new state-of-the-art performance in affordance grounding. It shows that our weighted context broadcasting module facilitate the framework to focus more on object-oriented information, effectively identifying the affordance regions.
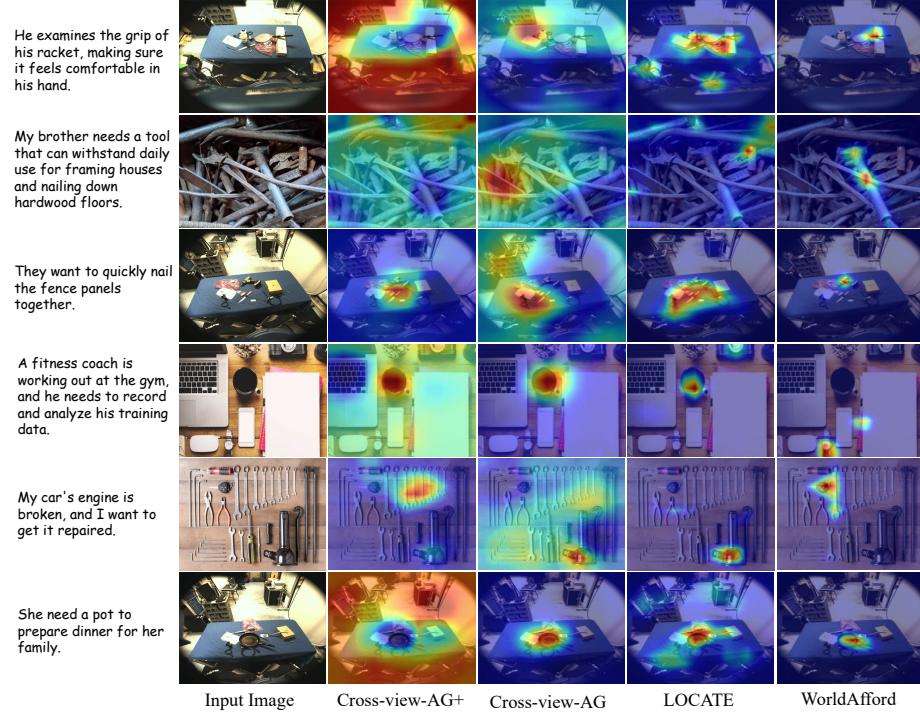
He examines the grip of his racket, making sure it feels comfortable in his hand.

My brother needs a tool that can withstand daily use for framing houses and nailing down hardwood floors.

They want to quickly nail the fence panels together.

A fitness coach is working out at the gym, and he needs to record and analyze his training data.

My car's engine is broken, and I want to get it repaired.

She need a pot to prepare dinner for her family.

| Input Image | Cross-view-AG+ | Cross-view-AG | LOCATE | WorldAfford |

**Fig. 8:** Visual comparison on the LLMaFF dataset. We manually assign labels to other methods since they cannot adopt the textual input. The labels, "swing", "carry", "catch", "pick up", "catch", and "carry" correspond to the first through sixth rows respectively. The previous approaches cannot predict the affordance regions based on the textual instructions, whereas our method successfully identifies the affordance regions corresponding to natural language instructions and performs outstandingly in the complex environments.

We further conduct the experiments on the new LLMaFF dataset to evaluate the performance of our method on the new task. Since the previous methods [22,32,33] cannot directly adopt the textual instruction as input, we manually select the appropriate action labels for comparison. The results shown in Tab. 2 demonstrate that our method can effectively localize the affordance regions of the objects in the complex scene images. While Cross-view-AG+ [33] achieves outstanding results on AGD20K, its performance on the LLMaFF dataset is less impressive. This indicates a decrease in its ability to accurately localize the affordance regions, as shown by its negative score (-0.194) on the NSS indicator. There are two explanations for the difference in the performance of Cross-view-AG+ between the two datasets: 1) The complexity of the new task, with numerous potential disruptions, presents a major challenge. 2) Cross-view-AG+ likely overfit the AGD20K dataset. The results in Tab. 2 show a significant decrease in the performance of Cross-view-AG [32] and LOCATE [22], which also demonstrate

that previous methods cannot localize the affordance regions of objects in the complex scene images.

## 5.4   Qualitative results

We show the qualitative comparisons on AGD20K in Fig. 7. The affordance regions identified by Cross-view-AG tend to be too large, sometimes including irrelevant regions. In contrast, LOCATE prefers to predict smaller regions but often fails to capture the full affordance region of objects. Cross-view-AG+ is able to identify the regions associated with the action label but not accurately. In contrast, out framework, WorldAfford, achieves the new state-of-the-art performance in this simple setting, providing a sharper and more accurate results. It demonstrates that the weighted context broadcasting module (WCB) allows the framework to focus more on informative objects, thus capturing more knowledge of objects and localizing more accurately.

The results on the LLMaFF dataset are shown in Fig. 8. We find that Cross-view-AG+ fails to identify affordance regions of multiple objects in the images, resulting in disordered color distribution, and thus cannot provide effective visual information to the agent. Cross-view-AG also shows some failure cases. It can capture information about objects in the image. However, we observe that the information tends to be biased towards objects with larger sample sizes in the training dataset. The results demonstrate the lack of comprehensive understanding of the objects in the image. LOCATE can capture the affordance of a few objects. However, it often activates the affordance regions of some irrelevant objects, and may interpret several objects as a single entity. In comparison, our method can predict the affordance maps that are more consistent with natural language instructions, more accurate, and is capable of localizing the affordance regions of multiple objects, thus providing richer visual information.

**Results on complex language instructions** We explore the challenge of affordance grounding based on complex language instructions that previous methods [33] cannot understand, as illustrated in Fig. 9. Unlike simple ones, these instructions require a deeper understanding of human knowledge, demonstrating the superior flexibility and creativity of our method. WorldAfford successfully identifies intricate affordance regions, and can highlight the complementarity of object interactions, such as using a knife and an apple together for slicing, which provides the embodied agent with detailed visual information that enhances its ability to follow complex instructions involving multiple sub-tasks. Our exploration of building a chair from wooden planks and nails illustrates how our method systematically identifies and activates the necessary affordance regions for sawing, measuring, and assembling, providing a comprehensive solution to multi-step tasks. This advance in affordance grounding may open new avenues for robotics and AI applications, significantly enriching the interaction between agents and their environment.

Resuls on complex language instructions



| | A man wants an apple and plans to slice it into pieces. | A carpenter wants to use wooden planks and nails to create a chair. |
|---|---|---|
| LLMs | He should first pick up the apple itself, and then use a sharp knife to cut the apple into pieces. And then cut the apple on a cutting board and prevent any injurey while minimizing mess. | He might need to carry a saw to cut them. And pick up the ruler to measure the length, width, and height of the chair parts, and then hit with hammer to drive the nails into the woolden planks and secure them together. |
| type | multi-objects complementary affordance. | multi-objects sequential affordance. |

**Fig. 9:** Affordance results based on difficult language instructions. While previous methods struggle to infer from difficult language instructions, our method demonstrates the capability to comprehend such instructions and accurately identify the affordance regions of multiple objects.

**Table 3:** Generalization ability comparison of WorldAfford with other state-of-the-art affordance grounding methods on AGD20K.

| Approach | Input Instruction | KLD↓ | SIM↑ | NSS↑ |
|---|---|---|---|---|
| Hotspots [37] | Action Label | 1.994 | 0.237 | 0.577 |
| Cross-view-AG [32] | Action Label | 1.787 | 0.285 | 0.829 |
| Affcorr [14] | Action Label | 1.618 | 0.348 | 1.021 |
| LOCATE [22] | Action Label | 1.405 | 0.372 | 1.157 |
| WorldAfford(ours) | Action Label | **1.393** | **0.38** | **1.225** |

## 5.5  Generalization ability and learnable parameters

To evaluate the generalization ability of our method, we add the results of the unseen test on AGD20K, which is shown in Tab. 3. Additionally, all LL-MaFF images, including various scenes and many object categories such as nail gun, smartwatch and so on, are unseen in training, which also demonstrates the the superior generalization ability of our method. We use the the knowledge of foundation models, the training cost is very low, and the comparison of learnable parameters: 120.03M(Cross-view-AG)/82.27M (Cross-view-AG+)/6.5M (LOCATE)/6.5M (WorldAfford).

## 5.6  Ablation Study

We conduct the ablative experiments on the LLMaFF dataset to validate the effectiveness of the affordance reasoning chain-of thought prompting(ARCoT).

The results shown in Tab. 4 demonstrate that the object information and the action information derived from the LLM via our affordance reasoning chain-of-thought prompting (ARCoT) can both improve the performance for the task of affordance grounding based on language instructions. We also validate that the proposed WCB module can enhance the perception of affordance regions by enabling the model to focus on more informative objects. Overall, our contributions significantly improve the affordance grounding capabilities of the model and establish a new state-of-the-art performance in the affordance grounding based on natural language instructions task. To verify our adjustments for masking off irrelevant objects, we conduct experiments on LLMaFF, the results is shown in Tab. 5.

**Table 4:** Ablation results of the proposed modules. LMA denotes the action information associated with the manipulated objects inferred from the LLM. WCB indicates the weighted context broadcasting module. LMO represents the object information inferred from the LLM.

| LMA | WCB | LMO | KLD↓ | SIM↑ | NSS↑ |
|-----|-----|-----|------|------|------|
|     |     |     | 3.073 | 0.105 | -0.059 |
|     | ✓   |     | 2.729 | 0.114 | 0.428 |
| ✓   |     |     | 2.768 | 0.124 | 0.303 |
| ✓   | ✓   |     | 2.335 | 0.155 | 0.981 |
|     | ✓   | ✓   | 2.336 | 0.180 | 1.081 |
| ✓   |     | ✓   | 1.700 | 0.256 | 2.325 |
| ✓   | ✓   | ✓   | **1.163** | **0.386** | **2.819** |

**Table 5:** The results of using entire images as input and masking off irrelevant objects on LLMaFF.

| Input | KLD↓ | SIM↑ | NSS↑ |
|-------|------|------|------|
| entire image | 2.752 | 0.134 | 0.27 |
| mask off | 1.163 | 0.386 | 2.819 |

## 6    Conclusion

In this paper, we introduce a new task of affordance grounding based on natural language instructions and propose a novel framework, WorldAfford. Our framework uses LLMs to process natural language instructions and employs SAM and

CLIP for object segmentation and selection. We further propose a Weighted Context Broadcasting module, allowing WorldAfford to localize affordance regions of multiple objects. Additionally, we present a new dataset, LLMaFF, to benchmark this task. The experimental results demonstrate that WorldAfford outperforms the other state-of-the-art methods for affordance grounding on both the AGD20K dataset and the new LLMaFF dataset.

# References

1. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023) 7, 11

2. Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., Finn, C., Fu, C., Gopalakrishnan, K., Hausman, K., et al.: Do as i can, not as i say: Grounding language in robotic affordances. arXiv preprint arXiv:2204.01691 (2022) 4

3. Bylinskii, Z., Judd, T., Oliva, A., Torralba, A., Durand, F.: What do different evaluation metrics tell us about saliency models? IEEE transactions on pattern analysis and machine intelligence **41**(3), 740–757 (2018) 11

4. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 9650–9660 (2021) 10, 11

5. Chen, J., Gao, D., Lin, K.Q., Shou, M.Z.: Affordance grounding from demonstration video to target image. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6799–6808 (2023) 4

6. Chuang, C.Y., Li, J., Torralba, A., Fidler, S.: Learning to act properly: Predicting and explaining affordances from images. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 975–983 (2018) 4

7. Chung, H.W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, Y., Wang, X., Dehghani, M., Brahma, S., et al.: Scaling instruction-finetuned language models. arXiv preprint arXiv:2210.11416 (2022) 4

8. Cong, Y., Liao, W., Rosenhahn, B., Yang, M.Y.: Learning similarity between scene graphs and images with transformers. arXiv preprint arXiv:2304.00590 (2023) 2

9. Deng, S., Xu, X., Wu, C., Chen, K., Jia, K.: 3d affordancenet: A benchmark for visual object affordance understanding. In: proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 1778–1787 (2021) 4

10. Do, T.T., Nguyen, A., Reid, I.: Affordancenet: An end-to-end deep learning approach for object affordance detection. In: 2018 IEEE international conference on robotics and automation (ICRA). pp. 5882–5889. IEEE (2018) 4

11. Fang, K., Wu, T.L., Yang, D., Savarese, S., Lim, J.J.: Demo2vec: Reasoning object affordances from online videos. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2139–2147 (2018) 4

12. Feng, G., Zhang, B., Gu, Y., Ye, H., He, D., Wang, L.: Towards revealing the mystery behind chain of thought: a theoretical perspective. Advances in Neural Information Processing Systems **36** (2024) 4

13. Grauman, K., Westbury, A., Byrne, E., Chavis, Z., Furnari, A., Girdhar, R., Hamburger, J., Jiang, H., Liu, M., Liu, X., et al.: Ego4d: Around the world in 3,000 hours of egocentric video. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18995–19012 (2022) 5

14. Hadjivelichkov, D., Zwane, S., Agapito, L., Deisenroth, M.P., Kanoulas, D.: One-shot transfer of affordance regions? affcorrs! In: Conference on Robot Learning. pp. 550–560. PMLR (2023) 5, 11, 15

15. Hassanin, M., Khan, S., Tahtali, M.: Visual affordance and function understanding: A survey. ACM Computing Surveys (CSUR) **54**(3), 1–35 (2021) 2, 4

16. Huang, S., Jiang, Z., Dong, H., Qiao, Y., Gao, P., Li, H.: Instruct2act: Mapping multi-modality instructions to robotic actions with large language model. arXiv preprint arXiv:2305.11176 (2023) 4

17. Huang, W., Abbeel, P., Pathak, D., Mordatch, I.: Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In: International Conference on Machine Learning. pp. 9118–9147. PMLR (2022) 4

18. Huang, W., Wang, C., Zhang, R., Li, Y., Wu, J., Fei-Fei, L.: Voxposer: Composable 3d value maps for robotic manipulation with language models. arXiv preprint arXiv:2307.05973 (2023) 4

19. Hyeon-Woo, N., Yu-Ji, K., Heo, B., Han, D., Oh, S.J., Oh, T.H.: Scratching visual transformer's back with uniform attention. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5807–5818 (2023) 10

20. Khan, M., Qiu, Y., Cong, Y., Abu-Khalaf, J., Suter, D., Rosenhahn, B.: Segment any object model (saom): Real-to-simulation fine-tuning strategy for multi-class multi-instance segmentation. arXiv preprint arXiv:2403.10780 (2024) 4

21. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. arXiv preprint arXiv:2304.02643 (2023) 3, 7, 9, 11

22. Li, G., Jampani, V., Sun, D., Sevilla-Lara, L.: Locate: Localize and transfer object parts for weakly supervised affordance grounding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10922–10931 (2023) 2, 4, 5, 9, 10, 11, 12, 13, 15

23. Li, G., Sun, D., Sevilla-Lara, L., Jampani, V.: One-shot open affordance learning with foundation models. arXiv preprint arXiv:2311.17776 (2023) 4

24. Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: International Conference on Machine Learning. pp. 12888–12900. PMLR (2022) 2

25. Li, Y.L., Xu, Y., Xu, X., Mao, X., Yao, Y., Liu, S., Lu, C.: Beyond object recognition: A new benchmark towards object concept learning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 20029–20040 (2023) 4

26. Li, Z., Peng, B., He, P., Galley, M., Gao, J., Yan, X.: Guiding large language models via directional stimulus prompting. Advances in Neural Information Processing Systems **36** (2024) 4

27. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. In: NeurIPS (2023) 2

28. Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Li, C., Yang, J., Su, H., Zhu, J., et al.: Grounding dino: Marrying dino with grounded pre-training for open-set object detection. arXiv preprint arXiv:2303.05499 (2023) 4

29. Lu, L., Zhai, W., Luo, H., Kang, Y., Cao, Y.: Phrase-based affordance detection via cyclic bilateral interaction. IEEE Transactions on Artificial Intelligence (2022) 4

30. Lüddecke, T., Ecker, A.: Image segmentation using text and image prompts. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 7086–7096 (2022) 4

31. Luo, H., Zhai, W., Zhang, J., Cao, Y., Tao, D.: One-shot affordance detection. arXiv preprint arXiv:2106.14747 (2021) 4

32. Luo, H., Zhai, W., Zhang, J., Cao, Y., Tao, D.: Learning affordance grounding from exocentric images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2252–2261 (2022) 2, 3, 4, 5, 10, 11, 12, 13, 15

33. Luo, H., Zhai, W., Zhang, J., Cao, Y., Tao, D.: Grounded affordance from exocentric view. International Journal of Computer Vision pp. 1–25 (2023) 2, 4, 5, 10, 11, 12, 13, 14

34. Luo, H., Zhai, W., Zhang, J., Cao, Y., Tao, D.: Learning visual affordance grounding from demonstration videos. IEEE Transactions on Neural Networks and Learning Systems (2023) 4

35. Mi, J., Liang, H., Katsakis, N., Tang, S., Li, Q., Zhang, C., Zhang, J.: Intention-related natural language grounding via object affordance detection and intention semantic extraction. Frontiers in Neurorobotics **14**, 26 (2020) 4

36. Mi, J., Tang, S., Deng, Z., Goerner, M., Zhang, J.: Object affordance based multimodal fusion for natural human-robot interaction. Cognitive Systems Research **54**, 128–137 (2019) 4

37. Nagarajan, T., Feichtenhofer, C., Grauman, K.: Grounded human-object interaction hotspots from video. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8688–8697 (2019) 4, 11, 15

38. Nguyen, A., Kanoulas, D., Caldwell, D., Tsagarakis, N.: Object-based affordances detection with convolutional neural networks and dense conditional random fields (09 2017). https://doi.org/10.1109/IROS.2017.8206484 5

39. Nguyen, A., Kanoulas, D., Caldwell, D.G., Tsagarakis, N.G.: Object-based affordances detection with convolutional neural networks and dense conditional random fields. In: 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 5908–5915. IEEE (2017) 4

40. Nguyen, T., Vu, M.N., Vuong, A., Nguyen, D., Vo, T., Le, N., Nguyen, A.: Open-vocabulary affordance detection in 3d point clouds. In: 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 5692–5698. IEEE (2023) 4

41. Peters, R.J., Iyer, A., Itti, L., Koch, C.: Components of bottom-up gaze allocation in natural images. Vision research **45**(18), 2397–2416 (2005) 11

42. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021) 2, 3, 7, 9, 11

43. Ren, T., Liu, S., Zeng, A., Lin, J., Li, K., Cao, H., Chen, J., Huang, X., Chen, Y., Yan, F., et al.: Grounded sam: Assembling open-world models for diverse visual tasks. arXiv preprint arXiv:2401.14159 (2024) 4

44. Swain, M.J., Ballard, D.H.: Color indexing. International journal of computer vision **7**(1), 11–32 (1991) 11

45. Van Vo, T., Vu, M.N., Huang, B., Nguyen, T., Le, N., Vo, T., Nguyen, A.: Open-vocabulary affordance detection using knowledge distillation and text-point correlation. arXiv preprint arXiv:2309.10932 (2023) 4

46. Wang, S., Zhou, Z., Kan, Z.: When transformer meets robotic grasping: Exploits context for efficient grasp detection. IEEE robotics and automation letters **7**(3), 8170–8177 (2022) 4

47. Wang, S., Zhou, Z., Li, B., Li, Z., Kan, Z.: Multi-modal interaction with transformers: bridging robots and human with natural language. Robotica **42**(2), 415–434 (2024) 4

48. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q.V., Zhou, D., et al.: Chain-of-thought prompting elicits reasoning in large language models. Advances in Neural Information Processing Systems **35**, 24824–24837 (2022) 4

49. Zhai, W., Luo, H., Zhang, J., Cao, Y., Tao, D.: One-shot object affordance detection in the wild. International Journal of Computer Vision **130**(10), 2472–2500 (2022) 2

50. Zhang, Z., Zhang, A., Li, M., Zhao, H., Karypis, G., Smola, A.: Multimodal chain-of-thought reasoning in language models. arXiv preprint arXiv:2302.00923 (2023) 2, 4

51. Zhao, X., Cao, Y., Kang, Y.: Object affordance detection with relationship-aware network. Neural Computing and Applications **32**(18), 14321–14333 (2020) 4

52. Zhou, Z., Wang, S., Chen, Z., Cai, M., Kan, Z.: A robotic visual grasping design: Rethinking convolution neural network with high-resolutions. arXiv preprint arXiv:2209.07459 (2022) 4

53. Zhou, Z., Wang, S., Chen, Z., Cai, M., Kan, Z.: A novel framework for improved grasping of thin and stacked objects. IEEE Transactions on Artificial Intelligence (2023) 4

54. Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M.: Minigpt-4: Enhancing vision-language understanding with advanced large language models. arXiv preprint arXiv:2304.10592 (2023) 2