

Learning Affordance Grounding from Exocentric Images

Hongchen Luo^{1†*} Wei Zhai^{1‡} Jing Zhang^{3†} Yang Cao^{1,4†} Dacheng Tao^{2,3}
¹ University of Science and Technology of China
³ JD Explore Academy ² The University of Sydney
⁴ Institute of Artificial Intelligence, Hefei Comprehensive National Science Center

{lhc12, wzhai056}@mail.ustc.edu.cn, jing.zhang1@sydney.edu.au,
 forrest@ustc.edu.cn, dacheng.tao@gmail.com

Abstract

Affordance grounding, a task to ground (i.e., localize) action possibility region in objects, which faces the challenge of establishing an explicit link with object parts due to the diversity of interactive affordance. Human has the ability that transform the various exocentric interactions to invariant egocentric affordance so as to counter the impact of interactive diversity. To empower an agent with such ability, this paper proposes a task of affordance grounding from exocentric view, i.e., given exocentric human-object interaction and egocentric object images, learning the affordance knowledge of the object and transferring it to the egocentric image using only the affordance label as supervision. To this end, we devise a cross-view knowledge transfer framework that extracts affordance-specific features from exocentric interactions and enhances the perception of affordance regions by preserving affordance correlation. Specifically, an Affordance Invariance Mining module is devised to extract specific clues by minimizing the intra-class differences originated from interaction habits in exocentric images. Besides, an Affordance Co-relation Preserving strategy is presented to perceive and localize affordance by aligning the co-relation matrix of predicted results between the two views. Particularly, an affordance grounding dataset named AGD20K is constructed by collecting and labeling over 20K images from 36 affordance categories. Experimental results demonstrate that our method outperforms the representative models in terms of objective metrics and visual quality. Code: github.com/lhc1224/Cross-View-AG.

1. Introduction

The goal of affordance grounding is to locate the region of “action possibilities” of an object. For an intelligent

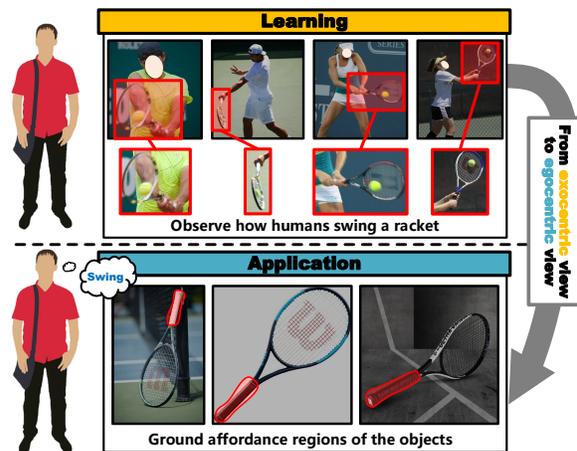


Figure 1. **Observation.** By observing the exocentric diverse interactions, the human learns affordance knowledge determined by the object’s intrinsic properties and transfer it to the egocentric view.

agent, it is necessary to know not only what the object is but also to understand how it can be used [11]. Perceiving and reasoning about possible interactions in local regions of objects is the key to the shift from passive perception systems to embodied intelligence systems that actively interact with and perceive their environment [1, 33, 34, 38]. It has a wide range of applications for robot grasping, scene understanding, action prediction [12, 13, 19, 23, 28, 31, 47, 49].

As affordance is a dynamic property closely related to the interaction between humans and environment [13], it is difficult to understand how to interact with objects and establish an explicit link between the objects’ intrinsic properties and affordances [29]. However, humans can easily perceive the object’s affordance region by observing exocentric human-object interactions, and give an egocentric definition. As shown in Fig. 1, although different persons hold the racket in different positions due to their individual habits, the human observer can perceive swingable regions

*This work was done during an internship at JD Explore Academy.

†Corresponding author. ‡ Equal contribution.

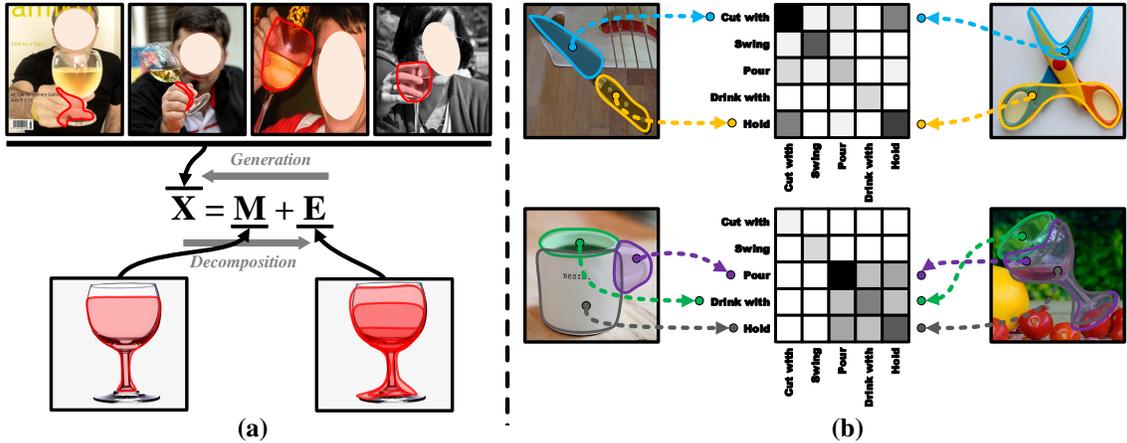


Figure 2. **Motivation.** (a) Exocentric interactions can be decomposed into affordance-specific features M and differences in individual habits E . (b) There are co-relations between affordances, *e.g.* “Cut with” inevitably accompanies “Hold” and is independent of the object category (knife and scissors). Such co-relation is common between objects. In this paper, we mainly consider extracting affordance-specific cues M from diverse interactions while preserving the affordance co-relations to enhance the perceptual capability of the network.

determined by the intrinsic properties (*e.g.*, the long handle structure) of the racket from a group of interacting images, despite the effect of individual differences, and transfer the knowledge to the egocentric view, thereby constructing a bridge between the object part and the affordance category.

To empower an agent with this ability to perceive the invariant egocentric affordance from various exocentric interactions, this paper proposes a task of affordance grounding from exocentric view, *i.e.*, given exocentric human-object interactions and egocentric object images, learning affordance knowledge and transferring it to object images by only using affordance labels as supervision. And in the testing stage, the output is the prediction of the affordance region for a specific object with the input of an egocentric object image and a particular affordance label.

To address this problem, we propose a cross-view knowledge transfer framework to extract affordance-specific features from exocentric interactions and transfer them to egocentric view. Specifically, we first devise an Affordance Invariance Mining (AIM) module to decompose the exocentric human-object interactions into the affordance representations determined by objects’ intrinsic properties and the differences originated from individual habits (as shown in Fig. 2 (a)). We use low rank matrix decomposition [10,18,22,24] to minimize the intra-class differences caused by diverse interactions to obtain affordance-specific cues. Furthermore, there is a correlation between the object affordances (as shown in Fig. 2 (b)), which can be adopted to establish the link between different affordances to reduce the uncertainty caused by multiple affordances regions on the object. Therefore, we present a novel Affordance Correlation Preserving (ACP) strategy to perceive and localize the affordance region by aligning the co-relation matrix of

prediction results from two views.

Despite the advances in affordance learning, the existing datasets [8, 29, 32, 35, 43] still bear limitations in terms of affordance/object category, image quality, and scene complexity. To carry out a comprehensive study, this paper proposes an affordance grounding dataset named AGD20K, consisting of 20,061 exocentric images and 3,755 egocentric images from 36 affordance categories. The contrastive experiments against several representative methods are performed on the AGD20K dataset. The results demonstrate the superiority of our proposed method in capturing the intrinsic property of objects and suppressing the interactive diversity of affordance.

Contributions: (1) We present a new affordance grounding from exocentric view task and establish a large-scale AGD20K benchmark to facilitate the research for empowering the agent to capture affordance knowledge from exocentric human-object interactions. (2) We propose a novel cross-view knowledge transfer framework for affordance grounding in which the affordance knowledge is acquired from exocentric human-object interactions and transferred to egocentric views while preserving the correlation between affordances, thereby achieving better perception and localization of interactive affordance. (3) Experiments on the AGD20K dataset demonstrate that our method outperforms state-of-the-art methods and can serve as a strong baseline for future research.

2. Related Works

2.1. Visual Affordance Grounding

The goal of affordance grounding is to locate the region of “action possibilities” of an object. Numerous works [4, 7, 8, 20, 27, 35, 48, 50] mainly build upon super-

vised approaches to establish mapping relations between local regions of objects and affordance. Sawatzky et al. [42, 43] adopt an Expectation-Maximization algorithm [6] to achieve weakly supervised affordance detection using only a few key points. Nagarajan et al. [33] exploit only affordance labels to ground the interactions from the videos. In contrast to [33], our goal is to empower the agent to learn affordance knowledge from exocentric human-object interactions. To this end, we propose an explicit cross-view knowledge transfer framework that extracts affordance knowledge determined by the intrinsic properties of objects from multiple exocentric interactions and transfers it into egocentric images.

2.2. Visual Affordance Dataset

The emergence of the relevant datasets drives the development of affordance grounding, as shown in Table 1. For example, Sawatzky et al. [43] select video frames from CAD120 [19] to construct a weakly supervised affordance detection dataset, using only cropped out object regions but in inferior image quality. Other affordance-related datasets [4, 8, 32, 35, 41] face the problems of small scale and low affordance/object category diversity and do not consider human actions to reason about the affordance regions. PAD dataset [29] considers the inference of human purpose from support images of human-object interactions and transfers to a group of query images but does not provide part-level affordance labels. In contrast to the above works, we explicitly consider exocentric-to-egocentric viewpoint transformations and collect a much larger scale of images, with richer affordance/object categories and part-level annotations, which are more useful and applicable to real-world application domains.

2.3. Learning View Transformations

The existing learning-view transformation works start from the theory of mirror neurons [40], which adopts embedding learning to generate perspective invariant representations from paired data, and leverage it for tasks such as action recognition and video summarization under egocentric view [16, 39, 44, 45]. For example, Li et al. [25] extract key egocentric signals from the exocentric view dataset during pre-training and distill them to the backbone to guide feature learning in the egocentric video task. In contrast to the above works, we aim to extract affordance knowledge from the diverse exocentric human-object interactions and transfer it to the egocentric view, which is challenging due to the uncertainty caused by various interactions and the multiple affordance regions that objects contain.

3. Method

Our goal is to ground the object affordance regions in egocentric images. During training, given a group of ex-

Table 1. **Statistics of related datasets and the proposed AGD20K dataset.** Part: part-level annotation. HQ: high-quality annotation. #Obj: number of object classes. #Aff: number of affordance classes. #Img: number of images.

Dataset	Year	HQ	Part	#Obj.	#Aff.	#Img.
UMD [32]	2015	✗	✓	17	7	30,000
[43]	2017	✗	✓	17	7	3,090
IIT-AFF [35]	2017	✗	✓	10	9	8,835
ADE-Aff [4]	2018	✓	✓	150	7	10,000
PAD [29]	2021	✓	✗	72	31	4,002
AGD20k (Ours)	2021	✓	✓	50	36	23,816

ocentric images $\mathcal{I}_{exo} = \{I_1, \dots, I_N\}$ (N is the number of exocentric images) and an egocentric object image I_{ego} , the network uses only affordance labels as supervision, so as to learn affordance knowledge from exocentric images and transfer it to egocentric images. During testing, only given an egocentric image I_{ego} and the affordance label C_a , the network outputs the affordance region on the object.

Our proposed cross-view knowledge transfer framework for affordance grounding is shown in Fig. 3. During training, we first use Resnet50 [14] to extract the features of exocentric and egocentric images to obtain $\mathcal{Z}_{exo} = \{Z_1, \dots, Z_N\}$ and Z_{ego} , respectively. We then present the Affordance Invariance Mining (AIM) module (see in Sec. 3.1) to extract affordance-specific clues (\mathcal{F}_{exo}) from the exocentric features. Meanwhile, we use two convolutional layers to map the egocentric feature to the embedding space consistent with the exocentric view: $F_{ego} = Conv(Z_{ego})$. Subsequently, the features of the two branches (\mathcal{F}_{exo} and F_{ego}) are fed into the same convolution layer to obtain features \mathcal{D}_{exo} and D_{ego} respectively. To ensure the affordance knowledge can be transferred to the egocentric view, we average the \mathcal{D}_{exo} through the global average pooling (GAP) layer to obtain the f_{exo} and pass the D_{ego} through the GAP layer to get the f_{ego} , and align f_{exo} and f_{ego} using L2 loss L_{KT} . Then, f_{exo} and f_{ego} are fed into the same fully connected layer to obtain the affordance prediction. Finally, we propose an Affordance Co-relation Preserving (ACP) strategy (see in Sec. 3.2) to enhance the network’s perception of affordance by aligning the co-relation matrix of the outputs of the two views. During testing, we feed the egocentric object images into the network only through the egocentric branch, and then use the CAM [51] technique to obtain the affordance regions of the object (see in Sec. 3.3).

3.1. Affordance Invariance Mining Module

As shown in Fig. 3, we decompose the interactions in exocentric images into affordance-specific features M and individual differences E . Inspired by low-rank matrix decomposition [10, 18, 22], we represent the M as the multiplication of a dictionary matrix W and a corresponding matrix H , where the dictionary bases represent the sub-features of human-object interaction, and minimize E by iterative op-

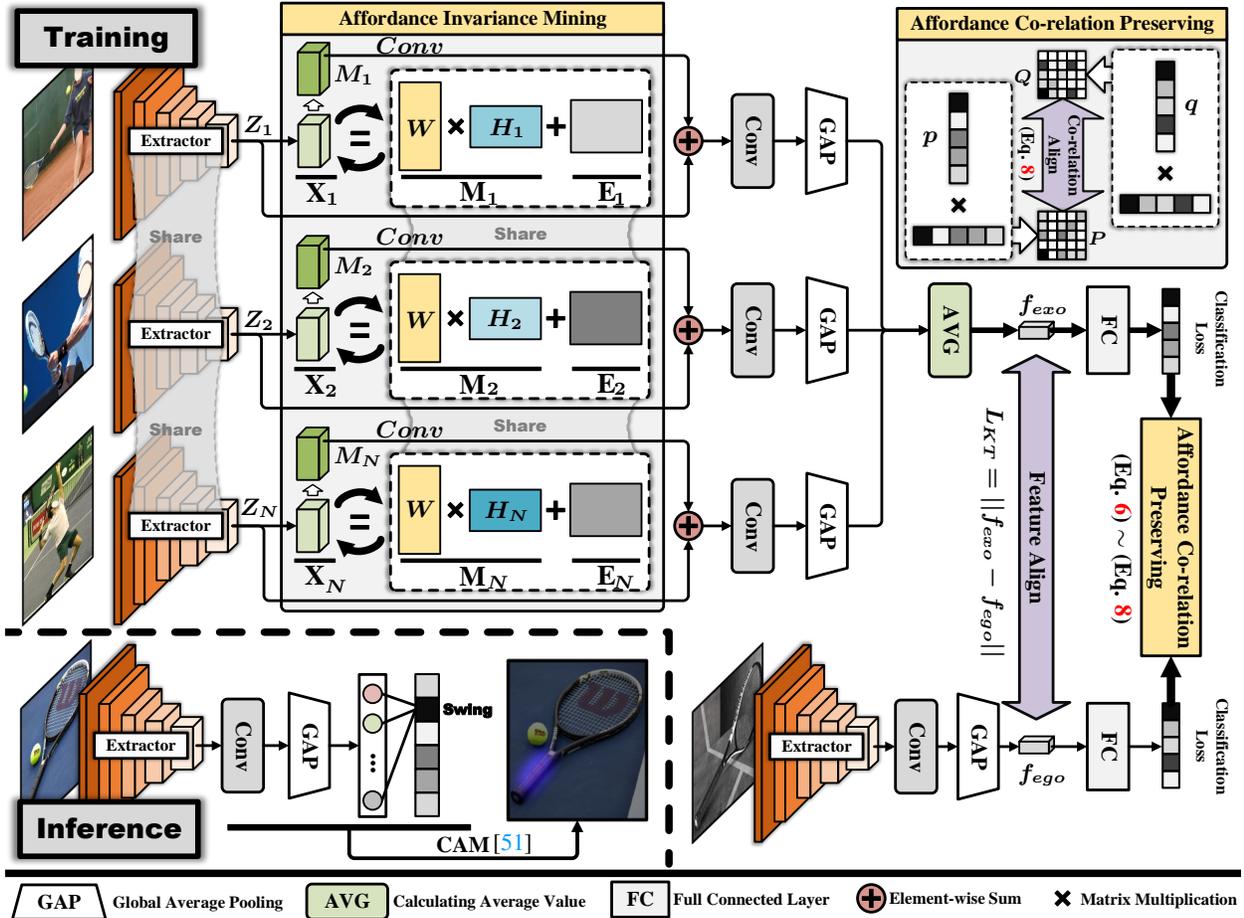


Figure 3. **Overview of the proposed cross-view knowledge transfer affordance grounding framework.** It mainly consists of an Affordance Invariance Mining (AIM) module and an Affordance Co-relation Preservation (ACP) strategy. The AIM module (see in Sec. 3.1) aims to obtain invariant affordance representations from diverse exocentric interactions. The ACP strategy (see in Sec. 3.2) enhances the network’s affordance perception by aligning the co-relation of the outputs of the two views.

timization to obtain a reconstructed affordance representation M . Specifically, for the input Z_i , we first reduce its dimensionality with a convolution layer and a ReLU layer to ensure the non-negativity of the input, and then reshape them into $X_i \in R^{c \times h \times w}$ (c , h and w are the channels, length, and width of the feature maps respectively). We use non-negative matrix factorization (NMF) [22] to update the dictionary and the coefficient matrices. Consequently, X_i is decomposed into two non-negative matrices W and H_i . Here $W \in R^{c \times r}$ is the dictionary matrix shared by all exocentric features, while $H_i \in R^{r \times h \times w}$ is the coefficient matrix of each exocentric feature, and r is the rank of the low-rank matrix W . To update H_i and W in parallel, we concatenate $\mathcal{X}_{exo} = \{X_1, \dots, X_N\}$ and $\mathcal{H} = \{H_1, \dots, H_N\}$ to obtain $X \in R^{c \times N \times h \times w}$ and $H \in R^{r \times N \times h \times w}$. Mathematically, the optimization process can be formulated as follows:

$$\min_{W, H} \|X - WH\|, \quad s.t. W_{ab} \geq 0, H_{bk} \geq 0. \quad (1)$$

W and H are updated according to the following rules:

$$H_{ab} \leftarrow H_{ab} \frac{(W^T X)_{ab}}{(W^T W H)_{ab}}, W_{ab} \leftarrow W_{ab} \frac{(X H^T)_{ab}}{(W H H^T)_{ab}}. \quad (2)$$

After several iterations, we get the output $M = WH$, and reshape it to $\mathcal{M}_{exo} = \{M_1, \dots, M_N\}$, $M_i \in R^{c \times h \times w}$. Finally, we use a convolution layer to map it to the residual space and sum it with the \mathcal{Z} to get the final output \mathcal{F}_{exo} :

$$F_i = Z_i + Conv(M_i), \quad i \in [1, N]. \quad (3)$$

In each batch of training, we update the initial dictionary matrix $W^{(0)}$ such that it can contain the statistical prior of the common subfeature of human-object interaction, *i.e.*,

$$W^{(0)} \leftarrow \alpha W^{(0)} + (1 - \alpha) \bar{W}, \quad (4)$$

where \bar{W} is the average over each mini-batch.

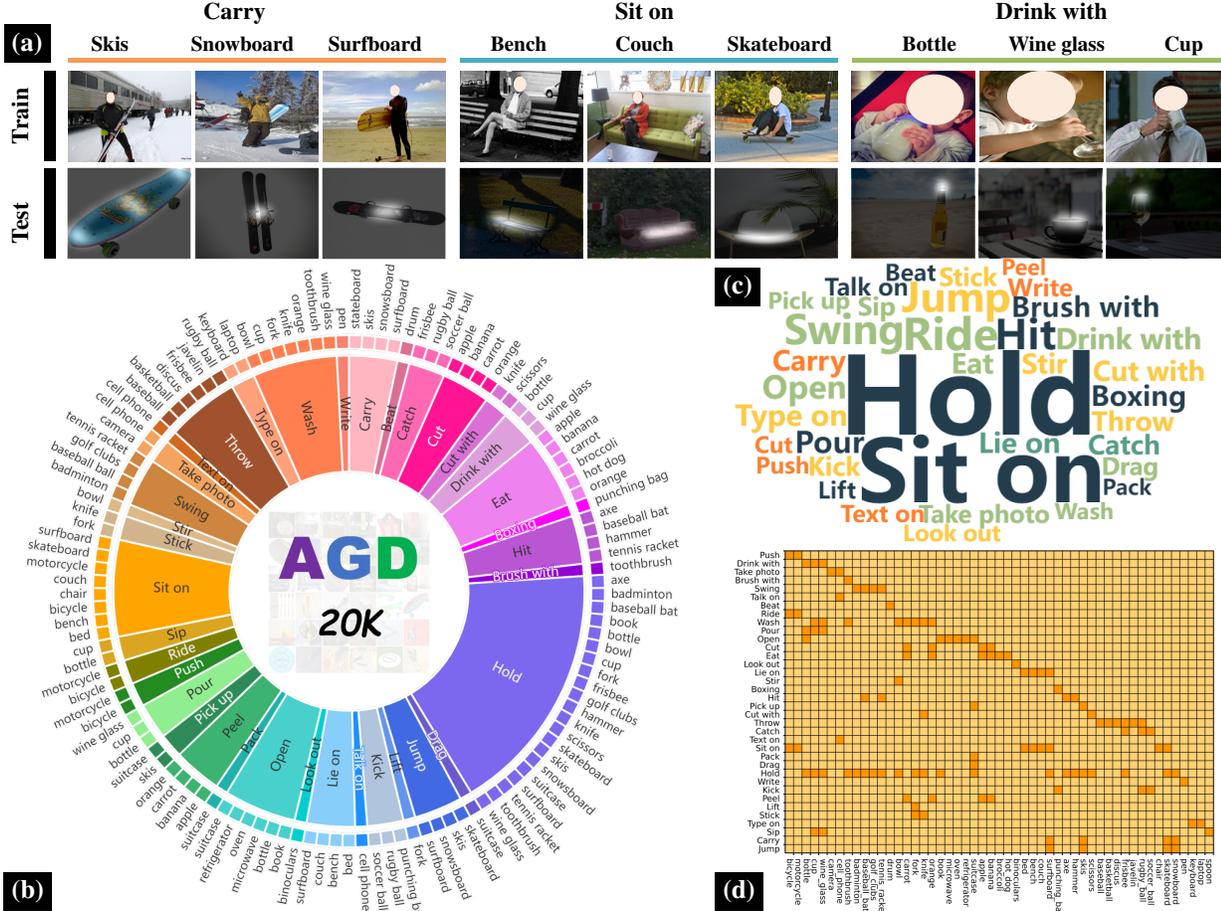


Figure 4. **The properties of the AGD20K dataset.** (a) Some examples from the dataset. (b) The distribution of categories in AGD20K. (c) The word cloud distribution of affordances in AGD20K. (d) Confusion matrix between the affordance category and the object category in AGD20K, where the horizontal axis denotes the object category and the vertical axis denotes the affordance category.

3.2. Affordance Co-relation Preserving Strategy

As shown in Fig. 3, we feed the feature representations of the two branches (f_{exo} and f_{ego}) into the same fully connected layer respectively to obtain the affordance category prediction scores s and g :

$$s = FC(f_{exo}), \quad g = FC(f_{ego}). \quad (5)$$

Then, we align the affordance co-relation between the exocentric and egocentric views by calculating the cross-entropy loss [15] L_{ACP} of the co-relation matrix of the prediction scores of the two branches:

$$p_j = \frac{\exp(s_j/T)}{\sum_k^{N_c} \exp(s_k/T)}, \quad q_j = \frac{\exp(g_j/T)}{\sum_k^{N_c} \exp(g_k/T)}, \quad (6)$$

$$P = pp^T, Q = qq^T, \quad (7)$$

$$L_{ACP} = - \sum_j \sum_k P_{jk} \log(Q_{jk}), \quad (8)$$

where T is used to control the degree of attention paid to the correlations between negative labels. P_{jk} and Q_{jk} denote the correlation between classes j and k in the prediction results. Finally, the total loss can be calculated as:

$$L = \lambda_1 L_{cls} + \lambda_2 L_{ACP} + \lambda_3 L_{KT}, \quad (9)$$

where λ_1 , λ_2 and λ_3 are hyper-parameters to balance the classification loss, ACP loss and L_{KT} loss. L_{cls} is the sum of the cross-entropy losses of the classification results of the two branches, and L_{KT} is loss of cross-view affordance knowledge transfer: $L_{KT} = \|f_{exo} - f_{ego}\|$.

3.3. Inference

Our test procedure only requires an egocentric object image and an affordance label as input to predict the affordance region. We utilize the class activation mapping [51] by computing a weighted sum of the feature maps D^i of the last convolutional layer to obtain the affordance region heatmap: $Y^{C_a} = \sum_i w_i^{C_a} D^i$, where C_a is the affordance

Table 2. **The results of different methods on AGD20k.** The best results are in **bold**. “Seen” means that the training set and the test set contain the same object categories, while “Unseen” means that the object categories in the training set and the test set do not overlap. The \diamond defines the relative improvement of our method over other methods. “Dark red”, “Orange” and “Purple” represent saliency detection, weakly supervised object localization and affordance grounding models, respectively.

Methods	Seen			Unseen		
	KLD ↓	SIM ↑	NSS ↑	KLD ↓	SIM ↑	NSS ↑
Minet [5]	5.197 $\diamond 70.4\%$	0.280 $\diamond 19.3\%$	0.596 $\diamond 55.5\%$	5.012 $\diamond 64.3\%$	0.263 $\diamond 8.4\%$	0.595 $\diamond 39.3\%$
DeepGazeII [21]	1.858 $\diamond 17.2\%$	0.280 $\diamond 19.3\%$	0.623 $\diamond 48.8\%$	1.990 $\diamond 10.2\%$	0.256 $\diamond 11.3\%$	0.597 $\diamond 38.9\%$
EgoGaze [17]	4.185 $\diamond 63.2\%$	0.227 $\diamond 47.1\%$	0.333 $\diamond 178\%$	4.285 $\diamond 58.3\%$	0.211 $\diamond 35.1\%$	0.350 $\diamond 137\%$
EIL [30]	1.931 $\diamond 20.4\%$	0.285 $\diamond 17.2\%$	0.522 $\diamond 77.6\%$	2.167 $\diamond 17.5\%$	0.227 $\diamond 25.6\%$	0.330 $\diamond 151\%$
SPA [36]	5.528 $\diamond 72.2\%$	0.221 $\diamond 51.1\%$	0.357 $\diamond 160.6\%$	7.425 $\diamond 75.9\%$	0.169 $\diamond 68.6\%$	0.262 $\diamond 216\%$
TS-CAM [9]	1.842 $\diamond 16.5\%$	0.260 $\diamond 28.5\%$	0.336 $\diamond 176\%$	2.104 $\diamond 15.1\%$	0.201 $\diamond 41.8\%$	0.151 $\diamond 449\%$
Hotspots [33]	1.773 $\diamond 13.3\%$	0.278 $\diamond 20.1\%$	0.615 $\diamond 50.7\%$	1.994 $\diamond 10.4\%$	0.237 $\diamond 20.3\%$	0.577 $\diamond 43.7\%$
Ours	1.538 ± 0.017	0.334 ± 0.001	0.927 ± 0.007	1.787 ± 0.017	0.285 ± 0.002	0.829 ± 0.014

Table 3. **Ablation study.** We investigate the influence of the AIM module, ACP strategy and L_{KT} on model performance.

	AIM	ACP	L_{KT}	KLD ↓	SIM ↑	NSS ↑
Seen				1.985	0.238	0.302
	✓			1.750	0.280	0.674
		✓		1.810	0.257	0.687
			✓	1.933	0.241	0.344
	✓	✓		1.749	0.286	0.735
	✓		✓	1.664	0.309	0.818
	✓	✓	✓	1.741	0.299	0.679
	✓	✓	✓	1.538	0.334	0.927
Unseen				2.059	0.228	0.445
	✓			1.933	0.261	0.682
		✓		1.920	0.250	0.666
			✓	1.967	0.265	0.622
	✓	✓		1.926	0.269	0.696
	✓		✓	1.916	0.272	0.679
	✓	✓	✓	1.922	0.267	0.640
	✓	✓	✓	1.787	0.285	0.829

class, D^i is the i -th layer feature map, and $w_i^{C_a}$ is the weight corresponding to the i -th neuron under the C_a category.

4. Dataset

Dataset Collection. The exocentric images are mainly obtained from HICO [3] and COCO [26]. We select images from the HICO dataset according to the verb category and the COCO dataset according to the object category. Then, we manually remove images with ambiguous interactions. To enrich the diversity of the dataset, we further collect 2, 112 exocentric images from free-license websites. Meanwhile, We collect 3, 755 egocentric images from the Internet with free use license according to object categories. Some examples are shown in Fig. 4 (a).

Dataset Annotation. We select 36 affordance classes commonly used in real-world application scenarios and assign labels to each image based on the interaction between human and object in each exocentric image. Given the object class contained in each affordance class, we assign affordance labels based on the object class in the egocen-

tric images. The testing process requires pixel-level labels to calculate objective metrics. As the annotation approach in [8], we take the form of points for regions of interaction, in which the dense points are for regions of frequent interaction and vice versa. Then, heatmaps of affordance regions can be obtained from the points as [8]. Some annotation examples are shown in Fig. 4 (a).

Statistic Analysis. To obtain deeper insights into our AGD20K dataset, we show its important features from the following aspects. The distribution of categories in the dataset is shown in Fig. 4 (b), which shows that the dataset contains a wide range of affordance/object categories in diverse scenarios. The affordance word cloud is shown in Fig. 4 (c). The confusion matrix of affordance and object categories is shown in Fig. 4 (d). It shows a multi-to-multi relationship between affordance and object categories, posing a significant challenge for the affordance grounding task. See supplementary materials for more details.

5. Experiments

5.1. Benchmark Setting

To provide a comprehensive evaluation, we choose three commonly used metrics **Kullback-Leibler Divergence (KLD)** [2], **SIMilarity (SIM)** [46] and **Normalized Scanpath Saliency (NSS)** [37], see supplementary material for details of each metric. Our model is implemented in PyTorch and trained with the SGD optimizer. The input images are randomly clipped from 256×256 to 224×224 with random horizontal flipping. We train the model for 35 epochs on a single NVIDIA 1080ti GPU with an initial learning rate of $1e-3$. The hyper-parameters λ_1 , λ_2 and λ_3 are set to 1, 0.5 and 0.5 respectively. The hyper-parameter T in the ACP is set to 1. The rank r of the dictionary matrix W and the number of iterations in the AIM are set to 64 and 6 respectively. The number of exocentric images N is set to 3. Besides, three saliency detection models (**Minet** [5], **DeepGazeII** [21], **EgoGaze** [17]), three weakly supervised object localization models (**EIL** [30], **SPA** [36], **TS-CAM**

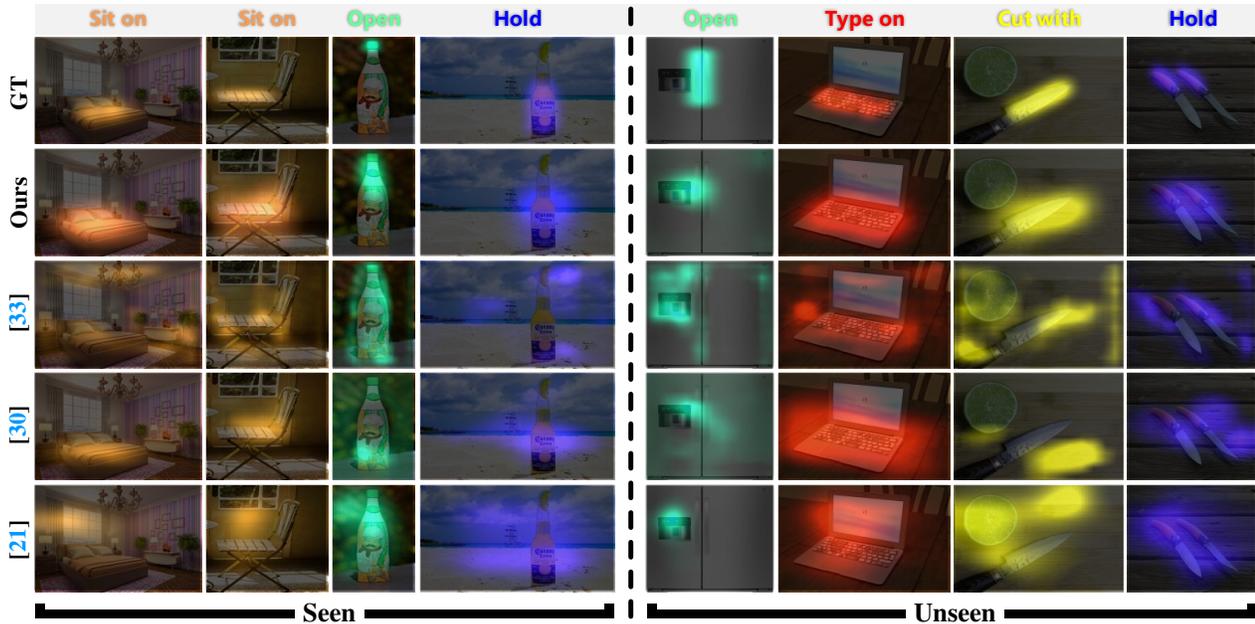


Figure 5. **Visual affordance heatmaps on the AGD20K dataset.** We select the prediction results of representative methods of affordance grounding (*Hotspots* [33]), weakly supervised object localization (*EIL* [30]), and saliency detection (*DeepGazeII* [21]) for presentation.

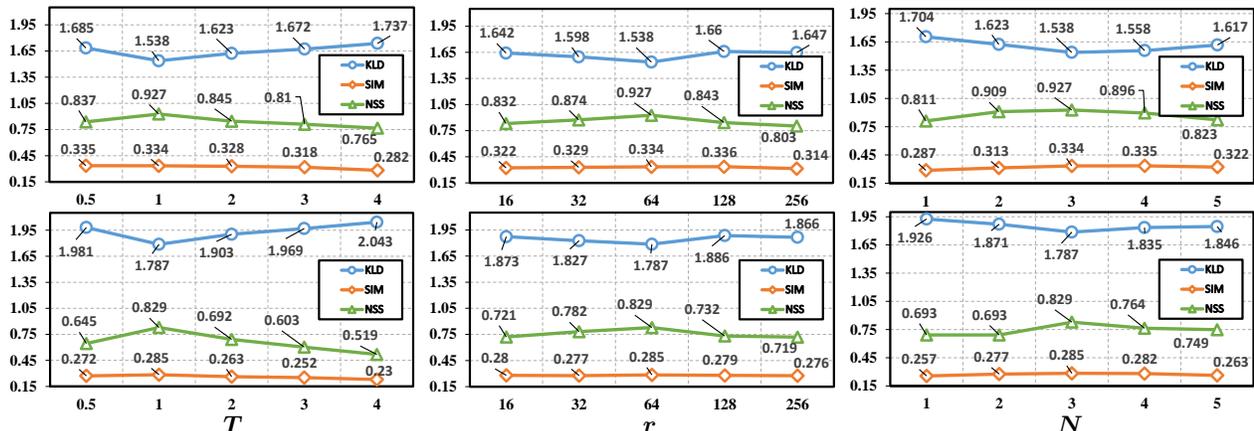


Figure 6. **Hyper-parameter study.** We investigate the influence of T in the ACP, the rank r of the W in the AIM, and the number of exocentric images N , respectively. The top and bottom columns represent the “Seen” and “Unseen” experimental settings, respectively.

[9]) and one affordance grounding model (*Hotspots* [33]) are chosen for comparison. We design two different settings: 1) “Seen”, *i.e.*, the training set and the test set contain the same object categories, and 2) “Unseen”, *i.e.*, the object categories in the training set and the test set do not overlap.

5.2. Quantitative and Qualitative Comparisons

The experimental results are shown in Table 2. Our method achieves the best results in both “Seen” and “Unseen” settings. Taking KLD as the metric, our method improves 17.2% compared to the best saliency model, 16.5% over the best weakly supervised object localization (WSOL) model, and 13.3% over the affordance grounding model in

the “Seen” setting. Our method with the “Unseen” setting improves 10.2% compared to the best saliency model, surpasses the best WSOL model by 15.1%, and exceeds the affordance grounding model by 10.4%. It indicates that our method can effectively transfer the affordance knowledge from the exocentric view to the object in the egocentric view and has a good generalization ability for unseen objects.

In addition, we visualize the affordance maps in “Seen” and “Unseen” settings, as shown in Fig. 5. It shows that our method can obtain more accurate prediction results for affordance grounding. While “Sit on” contains objects with different appearances (“bed” and “chair”), our method can capture the common features of the affordance region and

Table 4. **Different classes.** The KLD results of different methods on some representative affordance categories.

Classes	Hold	Swing	Drink with	Lie on	Brush with
Mlnet [5]	6.762	9.248	4.497	4.767	6.215
DeepGazeII [21]	2.071	2.478	2.067	1.602	2.385
EgoGaze [17]	4.671	6.723	4.268	2.921	5.135
EIL [30]	2.008	2.486	2.254	1.377	3.003
SPA [36]	3.006	6.720	7.683	4.006	8.043
TS-CAM [9]	1.628	2.420	2.300	1.370	2.642
Hotspots [33]	1.770	2.178	1.942	1.566	2.154
Ours	1.594	2.161	1.748	1.039	2.040

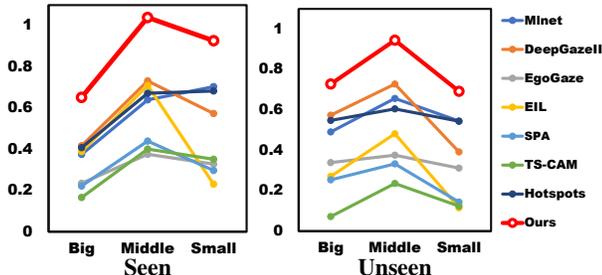


Figure 7. **Different scales.** We split the test set into three subsets of “Big”, “Middle” and “Small” according to the ratio of mask to the whole image, and show the results of the NSS metrics.

obtain better prediction results. Since “bottle” has two different affordances, “Open” and “Hold”, the network predicts different affordance regions. In the “Unseen” setting, the “knife” has two different affordances, “Hold” and “Cut with”. Our method can locate different affordance regions based on the learned affordance knowledge, demonstrating its superior generalization capability.

5.3. Ablation Study

The ablation study results are shown in Table 3. It confirms the ability of the AIM module to learn affordance-specific features from diverse exocentric interactions, which play a significant role in improving network performance. The ACP strategy improves more obviously than L_{KT} , indicating that the preservation of affordance co-relation can more effectively improve the network’s ability to perceive and locate affordance regions. In addition, we investigate the influence of different hyper-parameter settings of T in the ACP strategy (shown in Fig. 6 (left)), the rank r of the dictionary matrix W in the AIM module (shown in Fig. 6 (middle)), and different exocentric images N (see Fig. 6 (right)). It can be seen that the performance of the model is more influenced by T , while the rank r does not have a significant impact on the results. The number of exocentric images taken from 2 to 3 has a larger impact on the model. For $N = 1$, our model still outperforms most contenders.

Table 5. **Different sources.** “Exo” means that training only uses exocentric images, “Exo&Ego” means that training uses both exocentric and egocentric images.

	Method	Source	KLD ↓	SIM ↑	NSS ↑
Seen	EIL [30]	Exo	1.931	0.285	0.522
		Exo&Ego	2.156	0.321	0.747
	SPA [36]	Exo	5.528	0.221	0.357
		Exo&Ego	4.312	0.252	0.494
	TS-CAM [9]	Exo	1.842	0.260	0.336
		Exo&Ego	1.707	0.290	0.622
Ours	Exo&Ego	1.538	0.334	0.927	
Unseen	EIL [30]	Exo	2.167	0.227	0.330
		Exo&Ego	2.029	0.256	0.529
	SPA [36]	Exo	7.425	0.169	0.262
		Exo&Ego	6.174	0.209	0.433
	TS-CAM [9]	Exo	2.104	0.201	0.151
		Exo&Ego	2.002	0.228	0.305
Ours	Exo&Ego	1.787	0.285	0.829	

5.4. Performance Analysis

Different Classes. The KLD metrics on some representative categories are shown in Table 4. “Hold” and “Swing” both contain diverse object categories with different appearances. “Drink with” and “Hold” contain overlapped object categories but have completely different affordance regions. Objects of “Lie on” are generally labeled with a larger region, while those of “Brush with” are generally smaller. Our model exceeds others regarding different aspects of the challenge, which confirms its robustness. See supplementary material for the KLD metrics for each category.

Different Scales. We divide the test set into “big”, “middle” and “small” splits according to the proportion of mask to the whole image (see supplementary material for details). The test results are shown in Fig. 7. Our model outperforms all other methods in all splits on both settings, showing its ability to capture the intrinsic affordance properties of objects, even in more challenging cases. The performance of the experimental results on all metrics are shown in the supplementary material.

Different Sources. The results for different sources are shown in Table 5. It shows that using both exocentric and egocentric images improves most methods, but the improvement is limited. Our method still surpasses all models, showing that the knowledge transfer from explicitly cross-views is effective in learning from exocentric diverse interactions to egocentric invariant affordance representation. The performance of the experimental results on all metrics are shown in the supplementary material.

Limitations. Our method still has limitations, *e.g.*, the predicted affordance maps may contain intermediate background regions when multiple objects appear and irrelevant background regions may be activated for slender object. In

the future, we will refer to [36] to refine the generated results to obtain more accurate results.

6. Conclusion

In this paper, we make an attempt to address a new challenging task named affordance grounding from exocentric view. Specifically, we propose a novel cross-view knowledge transfer framework that can extract invariant affordance from diverse exocentric interactions and transfer it to egocentric view. We establish a large affordance grounding dataset named AGD20K, which contains 20K well-annotated images, serving as a pioneer testbed for the task. Our model outperforms representative models from related areas and can serve as a strong baseline for future research.

Broader Impacts. The research on affordance grounding from exocentric view will advance the realization of embodied intelligence. However, harmful human demonstrations (risky behaviors) may lead to negative guidance for the agent, which should be prohibited by strict legislation.

Acknowledgments. This work was supported by National Key R&D Program of China under Grant 2020AAA0105701, National Natural Science Foundation of China (NSFC) under Grants 61872327 and ARC FL-170100117.

References

- [1] Jeannette Bohg, Karol Hausman, Bharath Sankaran, Oliver Brock, Danica Kragic, Stefan Schaal, and Gaurav S Sukhatme. Interactive perception: Leveraging action in perception and perception in action. *IEEE Transactions on Robotics*, 33(6):1273–1291, 2017. 1
- [2] Zoya Bylinskii, Tilke Judd, Aude Oliva, Antonio Torralba, and Frédo Durand. What do different evaluation metrics tell us about saliency models? *IEEE transactions on pattern analysis and machine intelligence*, 41(3):740–757, 2018. 6
- [3] Yu-Wei Chao, Yunfan Liu, Xieyang Liu, Huayi Zeng, and Jia Deng. Learning to detect human-object interactions. In *2018 IEEE winter conference on applications of computer vision (WACV)*, pages 381–389. IEEE, 2018. 6
- [4] Ching-Yao Chuang, Jiaman Li, Antonio Torralba, and Sanja Fidler. Learning to act properly: Predicting and explaining affordances from images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 975–983, 2018. 2, 3
- [5] Marcella Cornia, Lorenzo Baraldi, Giuseppe Serra, and Rita Cucchiara. A deep multi-level network for saliency prediction. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 3488–3493. IEEE, 2016. 6, 8
- [6] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977. 3
- [7] Thanh-Toan Do, Anh Nguyen, and Ian Reid. Affordancenet: An end-to-end deep learning approach for object affordance detection. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 5882–5889. IEEE, 2018. 2
- [8] Kuan Fang, Te-Lin Wu, Daniel Yang, Silvio Savarese, and Lim J. Joseph. Demo2vec: Reasoning object affordances from online videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 3, 6
- [9] Wei Gao, Fang Wan, Xingjia Pan, Zhiliang Peng, Qi Tian, Zhenjun Han, Bolei Zhou, and Qixiang Ye. Ts-cam: Token semantic coupled attention map for weakly supervised object localization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2886–2895, October 2021. 6, 7, 8
- [10] Zhengyang Geng, Meng-Hao Guo, Hongxu Chen, Xia Li, Ke Wei, and Zhouchen Lin. Is attention better than matrix decomposition? *arXiv preprint arXiv:2109.04553*, 2021. 2, 3
- [11] James J Gibson. The theory of affordances. *Hilldale*, 1977. 1
- [12] Helmut Grabner, Juergen Gall, and Luc Van Gool. What makes a chair a chair? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1529–1536. IEEE, 2011. 1
- [13] Mohammed Hassanin, Salman Khan, and Murat Tahtali. Visual affordance and function understanding: A survey. *ACM Computing Surveys (CSUR)*, 54(3):1–35, 2021. 1
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 3
- [15] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 5
- [16] Hsuan-I Ho, Wei-Chen Chiu, and Yu-Chiang Frank Wang. Summarizing first-person videos from third persons’ points of view. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 70–85, 2018. 3
- [17] Yifei Huang, Minjie Cai, Zhenqiang Li, and Yoichi Sato. Predicting gaze in egocentric video by learning task-dependent attention transition. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 754–769, 2018. 6, 8
- [18] Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009. 2, 3
- [19] Hema Swetha Koppula, Rudhir Gupta, and Ashutosh Saxena. Learning human activities and object affordances from rgb-d videos. *The International Journal of Robotics Research*, 32(8):951–970, 2013. 1, 3
- [20] Hema S Koppula and Ashutosh Saxena. Physically grounded spatio-temporal object affordances. In *European Conference on Computer Vision (ECCV)*, pages 831–847. Springer, 2014. 2
- [21] Matthias Kümmerer, Thomas SA Wallis, and Matthias Bethge. Deepgaze ii: Reading fixations from deep features trained on object recognition. *arXiv preprint arXiv:1610.01563*, 2016. 6, 7, 8

- [22] Daniel D. Lee and H. Sebastian Seung. Algorithms for non-negative matrix factorization. In *NIPS*, 2000. 2, 3, 4
- [23] Xingming Li, Wei Zhai, and Yang Cao. A tri-attention enhanced graph convolutional network for skeleton-based action recognition. *IET Computer Vision*, 15(2):110–121, 2021. 1
- [24] Xia Li, Zhisheng Zhong, Jianlong Wu, Yibo Yang, Zhouchen Lin, and Hong Liu. Expectation-maximization attention networks for semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 9167–9176, 2019. 2
- [25] Yanghao Li, Tushar Nagarajan, Bo Xiong, and Kristen Grauman. Ego-exo: Transferring visual representations from third-person to first-person videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6943–6953, 2021. 3
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 740–755. Springer, 2014. 6
- [27] Liangsheng Lu, Wei Zhai, Hongchen Luo, Yu Kang, and Yang Cao. Phrase-based affordance detection via cyclic bilateral interaction. *arXiv preprint arXiv:2202.12076*, 2022. 2
- [28] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. Learning visual affordance grounding from demonstration videos. *arXiv preprint arXiv:2108.05675*, 2021. 1
- [29] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. One-shot object affordance detection. *arXiv preprint arXiv:2108.03658*, 2021. 1, 2, 3
- [30] Jinjie Mai, Meng Yang, and Wenfeng Luo. Erasing integrated learning: A simple yet effective approach for weakly supervised object localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8766–8775, 2020. 6, 7, 8
- [31] Priyanka Mandikal and Kristen Grauman. Learning dexterous grasping with object-centric visual affordances. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6169–6176. IEEE, 2021. 1
- [32] Austin Myers, Ching L Teo, Cornelia Fermüller, and Yiannis Aloimonos. Affordance detection of tool parts from geometric features. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1374–1381. IEEE, 2015. 2, 3
- [33] Tushar Nagarajan, Christoph Feichtenhofer, and Kristen Grauman. Grounded human-object interaction hotspots from video. In *Proceedings of the IEEE International Conference on Computer Vision (CVPR)*, pages 8688–8697, 2019. 1, 3, 6, 7, 8
- [34] Tushar Nagarajan and Kristen Grauman. Learning affordance landscapes for interaction exploration in 3d environments. *arXiv preprint arXiv:2008.09241*, 2020. 1
- [35] Anh Nguyen, Dimitrios Kanoulas, Darwin G Caldwell, and Nikos G Tsagarakis. 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)*, pages 5908–5915. IEEE, 2017. 2, 3
- [36] Xingjia Pan, Yingguo Gao, Zhiwen Lin, Fan Tang, Weiming Dong, Haolei Yuan, Feiyue Huang, and Changsheng Xu. Unveiling the potential of structure preserving for weakly supervised object localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11642–11651, 2021. 6, 8, 9
- [37] Robert J Peters, Asha Iyer, Laurent Itti, and Christof Koch. Components of bottom-up gaze allocation in natural images. *Vision research*, 45(18):2397–2416, 2005. 6
- [38] Santhosh K Ramakrishnan, Dinesh Jayaraman, and Kristen Grauman. An exploration of embodied visual exploration. *International Journal of Computer Vision*, 129(5):1616–1649, 2021. 1
- [39] Krishna Regmi and Mubarak Shah. Bridging the domain gap for ground-to-aerial image matching. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 470–479, 2019. 3
- [40] Giacomo Rizzolatti and Laila Craighero. The mirror-neuron system. *Annu. Rev. Neurosci.*, 27:169–192, 2004. 3
- [41] Anirban Roy and Sinisa Todorovic. A multi-scale cnn for affordance segmentation in rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 186–201. Springer, 2016. 3
- [42] Johann Sawatzky and Jurgen Gall. Adaptive binarization for weakly supervised affordance segmentation. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1383–1391, 2017. 3
- [43] Johann Sawatzky, Abhilash Srikantha, and Juergen Gall. Weakly supervised affordance detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 2, 3
- [44] Gunnar A Sigurdsson, Abhinav Gupta, Cordelia Schmid, Ali Farhadi, and Karteek Alahari. Actor and observer: Joint modeling of first and third-person videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7396–7404, 2018. 3
- [45] Bilge Soran, Ali Farhadi, and Linda Shapiro. Action recognition in the presence of one egocentric and multiple static cameras. In *Asian Conference on Computer Vision*, pages 178–193. Springer, 2014. 3
- [46] Michael J Swain and Dana H Ballard. Color indexing. *International Journal of Computer Vision (IJCV)*, 7(1):11–32, 1991. 6
- [47] Yuxiang Yang, Zhihao Ni, Mingyu Gao, Jing Zhang, and Dacheng Tao. Collaborative pushing and grasping of tightly stacked objects via deep reinforcement learning. *IEEE/CAA Journal of Automatica Sinica*, 9(1):135–145, 2021. 1
- [48] Wei Zhai, Hongchen Luo, Jing Zhang, Yang Cao, and Dacheng Tao. One-shot object affordance detection in the wild. *arXiv preprint arXiv:2108.03658*, 2021. 2
- [49] Jing Zhang and Dacheng Tao. Empowering things with intelligence: a survey of the progress, challenges, and opportunities in artificial intelligence of things. *IEEE Internet of Things Journal*, 8(10):7789–7817, 2020. 1

- [50] Xue Zhao, Yang Cao, and Yu Kang. Object affordance detection with relationship-aware network. *Neural Computing and Applications*, 32(18):14321–14333, 2020. [2](#)
- [51] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2921–2929, 2016. [3](#), [4](#), [5](#)