

# Contrastively Learning Visual Attention as Affordance Cues from Demonstrations for Robotic Grasping

Yantian Zha, Siddhant Bhambrri and Lin Guan

**Abstract**—Conventional works that learn grasping affordance from demonstrations need to explicitly predict grasping configurations, such as gripper approaching angles or grasping preshapes. Classic motion planners could then sample trajectories by using such predicted configurations. In this work, our goal is instead to fill the gap between affordance discovery and affordance-based policy learning by integrating the two objectives in an end-to-end imitation learning framework based on deep neural networks. From a psychological perspective, there is a close association between attention and affordance. Therefore, with an end-to-end neural network, we propose to learn affordance cues as visual attention that serves as a useful indicating signal of how a demonstrator accomplishes tasks, instead of explicitly modeling affordances. To achieve this, we propose a contrastive learning framework that consists of a Siamese encoder and a trajectory decoder. We further introduce a coupled triplet loss to encourage the discovered affordance cues to be more affordance-relevant. Our experimental results demonstrate that our model with the coupled triplet loss achieves the highest grasping success rate in a simulated robot environment. Our project website can be accessed at <sup>1</sup>.

## I. INTRODUCTION

Humans tend to understand objects and their parts from potentially applicable actions or motion primitives that can achieve effects for accomplishing a task. This phenomenon is abstracted as an ecological psychology concept called affordance established by J. J. Gibson ([9]). An affordance defines a mapping from an object feature to all applicable actions ([23]). Essentially, an affordance represents an object-action-effect relationship, which is an interactive procedure between an actor (e.g. a hand) and an object (e.g. a mug). Consider the example shown in Fig. 1, in a mug grasping task, a human teacher’s affordance biases (affordance-effect judgments) might vary with the shape and size of the mug – the mug-A is graspable from its handle, whereas the mug-B is graspable from its body. When a robot learns from human demonstrations, it would be beneficial if the robot also discovers such affordance bias behind human demonstrations and generate (visual) affordance cues to support its learning.

There are several works ([6], [31], [7]) that propose to learn affordance knowledge from human demonstrations

<sup>\*</sup>This research is supported in part by ONR grants N00014-19-1-2119, N00014-16-1-2892, N00014-18-1-2442, N00014-18-1-2840, N00014-9-1-2119, AFOSR grant FA9550-18-1-0067, and a JP Morgan AI Faculty Research grant. The first author thanks his advisor Prof. Subbarao Kambhampati for his encouragement and support, Prof. Heni Ben Amor for a helpful discussion, and Allen Z. Ren for the communications with the first author that speeds up reproducing the baseline work.

All authors are with the School of CS & AI, Arizona State University, United States of America. Email: {yzha3, sbhambr1, lguan9}@asu.edu

<sup>1</sup><https://sites.google.com/asu.edu/affordance-aware-imitation/project>

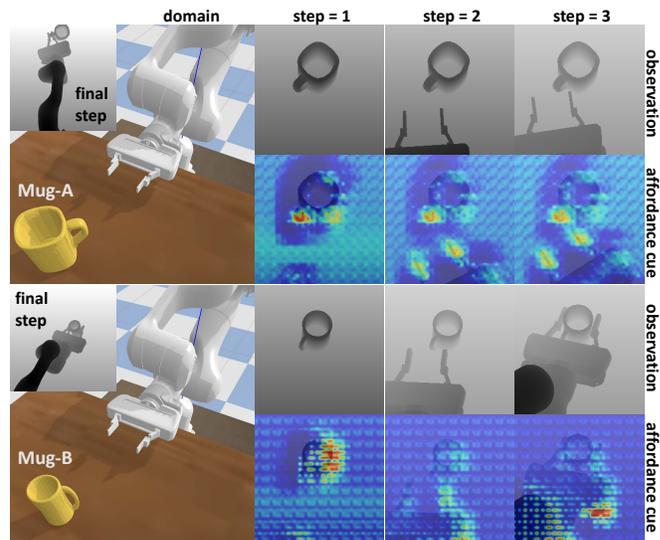


Fig. 1. An example of two mugs: Humans should have different affordance-effect judgments on their bodies or handles. At each step, the robot takes one depth image (observation) and generates one attention map, a.k.a affordance cue, that further triggers an action. We also show how the robot eventually picks up a mug, either by holding the body or handle, by showing the final step observation.

that avoid a labor-intensive process of collecting affordance labels. However, previous works fully decouple affordance discovery from behavior learning and execution, which means the affordance predictor and the robot controller are trained or constructed separately. In those methods, the predicted affordances are fed into classic motion planners in an engineered fashion. In this work, we instead combine the learning of affordance knowledge and motion generation from human demonstrations in an end-to-end deep imitation learning framework. We further argue that learning visual attention as affordance cues, rather than explicitly modeling affordances, is enough to be a reliable latent feature for the actor. The actor here is also a part of the whole neural network, instead of being a separate and unlearnable external planner.

When learning affordance knowledge from demonstrations, visual attention can be an informative cue for inferring affordance. In fact, the close association between visual attention and affordance has been investigated by some cognitive psychology works ([2], [14]). According to them, humans use visual attention conditioned on objects’ geometric and spatial properties to speed up affordance-effect judgments, which then helps generate motor signals (behaviors). When the attention comes from a controlled

mental process (i.e. top-down attention) driven by a task, object parts that are permissible for accomplishing the task would be highlighted, from which proper affordance-effect judgments can be derived ([26]). Indeed, humans do not explicitly think about what are all possible ways of picking-up a mug at each of its parts or pixels once she or he already learned how to grasp different mugs. An earlier work ([19]) also investigates how visual attention could be associated with affordance cueing in the context of robotics.

In this work, we hypothesize that by encouraging the robot to attend to discriminative features that explain the differences between different demonstrated behaviors, the robot will be able to more effectively discover affordance information and better imitate human behavior. This is because when a human teaches affordance knowledge to a robot, the human tends to think about what makes he/she have different affordance-effect judgments and how could the resulting trajectories be accordingly different. Again taking Fig. 1 as an example, if the human shows two different trajectories for the robot to pick up the two different mugs, the robot could discover important discriminative features regarding mug shapes and sizes and consequently better learn from demonstrations. If the robot could attend to the handle of Mug-A, then it would be more likely for the robot to grasp the Mug-A by its handle, rather than its body.

Therefore, our central problem is to help robots to imitate humans not only at the behavior level but also with a hidden objective of discovering the affordance-relevant distinctiveness underlying human demonstrations. As such, robot could attend to appropriate parts of a specific mug that hints on the same affordance in human’s mind and therefore helps trigger a similar behavior to humans’. To address the problem, we propose to use a deep Siamese encoder and trajectory decoder that are trained jointly with a contrastive loss and a behavior cloning loss in an end-to-end fashion. We also propose a coupled triplet loss to encourage the discovered discriminative features to be more affordance-relevant.

To thoroughly evaluate our work, we compare our model with a variant version that uses a normal triplet loss, a version without using Siamese network, and a baseline ConvNet-based behavior cloning model. We evaluate those models in terms of the grasping success rates and visualizations of predicted affordance cues. We empirically show that our model with the coupled triplet loss performs the best.

To the best of our knowledge, our work is the first to combine grasping affordance learning and imitation learning from expert demonstrations based on deep neural networks.

## II. RELATED WORKS

### A. Affordance Learning for Grasping

Regarding learning affordances for grasping, the majority of previous works use ground-truth affordance labels to learn affordances for grasping ([18], [34], [35], [30]). [18] use thermal maps to learn the graspable positions of several household objects. [34] employ transfer learning from pre-trained vision models to pixel-wise affordance prediction networks that help the robot generalize over novel objects.

[35] also uses pixel-level affordances to identify multi-grasp possibilities for objects present in a cluttered area.

There are also works that propose to learn affordances from demonstrations or via imitation learning: [10], [6], [31], [7]. All of these works share similar frameworks of using classic unsupervised learning (e.g. clustering) to identify gripper control parameters that would be fed into a motion planner. In [10], the affordance is represented by contact points for grasping an object. They use a predefined tracking configuration to reduce the number of potential contact points from demonstrations. After detecting a set of contact points on a new object, nearest-neighbor classification is used to identify a template grasp that matches their demonstration data on similar objects. In [6], [31], [7], they define affordable actions in affordances as approaching angles ([6]), grasp preshape hypotheses ([31]), or grasping prototypes ([7]). Then they use clustering to find condensed representations of affordances. In this work, we leverage the expressiveness of deep neural networks to implicitly learn affordance knowledge from human demonstrations. By learning from demonstrations with a deep contrastive learning framework, we evade the need of using ground truth affordances.

### B. Attention Guided Imitation Learning

Imitation learning or learning from demonstrations [20], [25] has been at the core of teaching robots to perform object-manipulation tasks in a similar way to humans performing the same task. By combining visual attention with imitation learning, robots could learn the information better by focusing on smaller but more important regions in manipulation tasks. In [1], authors use natural language descriptors that are specific to the task at hand to generate guided attention cues. The masked attention method places attention on the entire object that is to be grasped but does not focus on the specific region that the robot needs to interact with. A similar issue can be observed in [24] where the model first captures attention features of the object by generating attention maps for different stages of the object manipulation task. But the robot, in this case, can only learn to imitate the task and can not decide how to perform a specific task in a variable environment setting. Hence, one potential advantage of our approach is that it allows the robot to learn the specific graspable points in scenarios where the shape of the object (mug and its handle, in this case) can also vary restricting the possible graspable areas even for a human.

### C. Discriminative Feature Learning

The objective of discriminative feature learning is to make sure that the learned features of deep neural networks can represent different inputs contrastively enough [32]. Usually, such learned features can be easily separated by k-nearest neighbors algorithms [8]. Various approaches have been proposed to address the discriminative learning problem for deep neural networks. One popular approach is Siamese neural network [4] that was proposed in 1994 for verifying signatures. Quite a few works used Siamese neural networks as a backbone for new deep network models of discriminative

feature learning. Siamese neural network can be combined with ConvNets and trained with a Binary Cross Entropy loss as in [13], or triplet loss as in [27]). Recently deep learning based Siamese neural networks are also applied in many new applications, like face recognition ([27], [28]), object discovery ([29], [11]), object co-segmentation ([3], [17], [21]), and re-identification as in [36], [22].

Besides using Siamese neural networks for discriminative feature learning, [32] proposes to replace normal classification loss functions in ConvNets with the Center loss. Center loss works by minimizing the intra-class variations and meanwhile keeping the inter-class feature variations separable enough. The work [33] proposes a prototype-based discriminative feature learning (PDFL) method. The work [15] follows a similar idea to [32] and designs a discriminative feature learning algorithm for domain adaption.

### III. PROBLEM STATEMENT

Our goal is to encourage robots to learn from expert demonstrations more efficiently by taking advantage of discovering affordance-relevant distinctiveness underlying all demonstrations. To capture such distinctiveness, we could learn an attention model that predicts **affordance cues** from observations. We assume that humans have hidden affordance knowledge that is associated with visual cues. When focusing on a region of an object, humans also know what could be affordable grasps on that region. **The highlighting of such a region could serve as an affordance cue that supports the imitation learning of the learner itself.**

Our task is to pick up mugs in a way that allows pouring water in the near future. Hence, a robot could grasp a mug by reaching its gripper horizontally to the mug body, grasp the left and right sides of a handle, or grasp the front and back sides of a handle. Therefore, we have three candidate affordances (an object part and an applicable grasp): body-grasp, handle-left-right-grasp, and handle-front-back-grasp, as shown in Fig. 2.

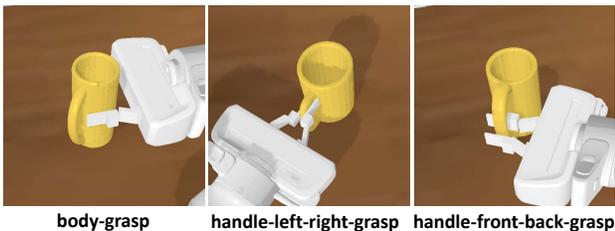


Fig. 2. The picture depicts three examples of the three candidate affordable grasps: body-grasp, handle-left-right-grasp, and handle-front-back-grasp.

We assume that robots share the same embodiment with humans: a human expert teleoperates a robot to collect demonstration trajectories. The collected trajectories are categorized in terms of the three candidate affordances. Each trajectory  $\tau$  is a sequence of triplets  $(o_t, s_t, a_t^*)$ :  $\tau := \{(o_t, s_t, a_t^*)\}_{t=1}^{T-1} \cup (s_T, o_T^*)$ .  $T$  denotes the trajectory length.  $o_t$  denotes a depth image at step  $t$ .  $s_t$  denotes a state vector at  $t$  which has eight values: the relative 3-D position of the

gripper to a mug, the relative 3-D Euler orientation of the gripper to mug, and two finger position values.  $a_t^*$  denotes an expert action vector that has seven values: the translation of  $x$ ,  $y$ , and  $z$ ; the rotation of roll, pitch, and yaw; and a value that indicates if the fingers should be closed or not. All trajectories are categorized into  $C$  **affordance categories** ( $C$  is three in this work). In each category  $c$ , there are  $N_c$  expert trajectories  $\{\tau_i^c\}_{i=1}^{N_c}$ .

### IV. APPROACH

Our framework learns to reproduce humans' behavior with visual cues that hint at different affordances from human demonstrations. The learning of such visual cues is achieved by training a Siamese encoder, and the policy imitation is by a behavior-cloning-based trajectory decoder. The Siamese encoder and trajectory decoder are trained simultaneously in a contrastive learning framework.

#### A. A Challenge and A Graphical Model

Since human embeds her/his hidden knowledge of affordances into trajectories of different affordance categories, intuitively we could use contrastive learning to help discover affordance cues. A well-known contrastive loss is the triplet loss ([27]) defined in Equation 1.

$$L(\mathcal{A}, \mathcal{P}, \mathcal{N}) = \sum_{i=1}^N [\|f(\mathcal{A}_i) - f(\mathcal{P}_i)\|_2^2 - \|f(\mathcal{A}_i) - f(\mathcal{N}_i)\|_2^2 + M]_+ \quad (1)$$

where  $\mathcal{A}$ ,  $\mathcal{P}$ , and  $\mathcal{N}$  denote the sets of anchor, positive, and negative trajectories;  $\mathcal{A}_i$  denotes the  $i$ -th trajectory in  $\mathcal{A}$  (likewise for  $\mathcal{P}_i$  and  $\mathcal{N}_i$ ); The  $i$ -th positive trajectory is sampled from the same affordance category of the  $i$ -th anchor trajectory, whereas the  $i$ -th negative trajectory is sampled from a different category.  $M$  denotes a margin value,  $\|z\|_2^2$  denotes a squared Euclidean distance metric,  $[z]_+$  denotes any  $z$  that is larger than zero, and,  $f()$  is an encoding function that can be parameterized by a neural network.

Contrastive losses encourage a learner to discover recurring patterns in one category and discriminative patterns across different categories. From humans' perspective, shifting affordance cues would lead to the change of affordance-effect judgment such that a different action would be taken. Thus, for two trajectories that come from the same affordance category, their affordance cues would share certain similarities; for two trajectories that are sampled from different affordance categories, their affordance cues tend to be different.

However, one potential challenge in using traditional contrastive learning frameworks is that the agent might learn to exploit affordance-irrelevant information (e.g. contexts, initial configurations) to distinguish two trajectories. Inspired by earlier affordance learning from demonstration works which extract grasping preshapes ([7]) or gripper-object approaching orientations ([6]) to learn affordance clusters, we also extract **the segment of interaction state-action pairs** (shortly **interaction segment**) of each trajectory as

an affordance-relevant feature. An interaction segment of a trajectory includes the state-action transitions that an actor (e.g. a gripper) changes the state of a target object. Since affordance is about object-action-effect relationships ([16]), we use the current state and previous expert action as the state-action features per step. Specifically speaking, we extract the interaction state-action transitions from the trajectory  $\tau$ :  $\{(s_t, a_{t-1}^*)\}_{t=m}^n$ . The curly brackets  $\{\}_{t=m}^n$  mean that the state of an object changes between the step  $m$  and  $n$  due to the motion of an actor.

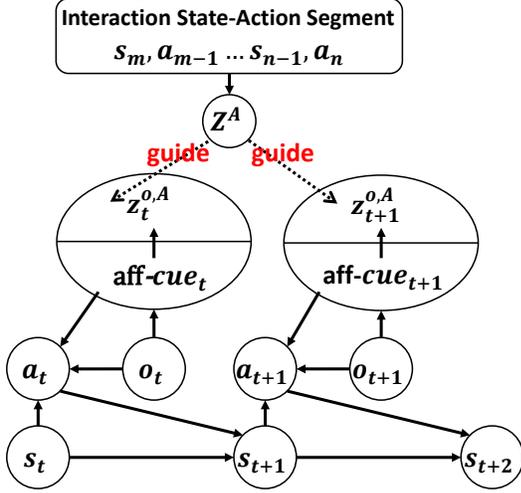


Fig. 3. The graphical model of Siamese encoder. Each  $\bigcirc$  denotes a node that represents a feature. Each directed edge from node  $X$  to node  $Y$  represents that “ $Y$ ” depends on “ $X$ ”. The  $\bigoplus$  denotes two nodes (features) that are generated by the same component one by one. The dashed edge means that the learning of affordance embedding  $Z^A$  guides the learning of observation embedding  $z_i^{o,A}$  at each step  $t$ .

By contrastively learning from such interaction segments, we could learn a clean representation of affordances in the form of embeddings. Such embeddings could be named **affordance embeddings**  $Z^A$  which is computed by using Equ. 2. We do not directly condition robots’ decision-making on  $Z^A$ . Instead, we treat  $Z^A$  as a guidance to help robots learn an attention model that extracts an affordance-cue from an observation for decision-making. To achieve this objective, we feed state  $s_t$ , image  $o_t$ , and action  $a_{t-1}^*$  at an arbitrary step  $t$  of a trajectory into a neural network that generates latent features: an affordance-cue and a processed visual feature. We then convert such latent features to **observation embedding**  $z_i^{o,A}$  as illustrated in Equ. 3. In the rest of this paper, we use  $Z^A$  and  $Z_\tau^A$  to denote the same thing: an affordance embedding for an arbitrary trajectory ( $\tau$ ). We use  $Z_\tau^A$  when we need distinguish among different trajectories. Likewise for  $z_i^{o,A}$  and  $z_{i,\tau}^{o,A}$ . Now we can encourage  $z_i^{o,A}$  to be either closer or farther from  $Z^A$  depending on if the trajectory that gives  $z_i^{o,A}$  belongs to the same affordance category of the trajectory that gives  $Z^A$  or not. This is essentially using the spatial relationships among embeddings in the space of  $Z^A$  to guide the learning of observation encoding (Equ. 3). This way, we link the visual processing of high-dimensional sensory at any step of a trajectory to

the affordance knowledge that can be learned faster from the lower-dimensional data of interaction segments.

$$Z_\tau^A = f^A(\text{interaction-segment}_\tau) = f^A(\{(s_t, a_{t-1}^*)\}_{t=m}^n | \tau) \quad (2)$$

$$z_{i,\tau}^{o,A} = f^O(s_t^\tau, o_t^\tau, a_{t-1}^{\tau,*}) \quad (3)$$

where  $f^A()$  and  $f^O()$  are two encoding functions whose neural network architectures are explained in Sec. IV-C;  $Z_\tau^A$  encodes the interaction segment of the trajectory  $\tau$ ; and  $z_{i,\tau}^{o,A}$  mainly encodes the high-dimension observation  $o_t$  (e.g. a depth image) at the current step  $t$  of the trajectory  $\tau$  with auxiliary information like the current state  $s_t$  and previous ground-truth action  $a_{t-1}^*$ .

The graphical model that describes the above idea is illustrated in Fig. 3 which shows two transitions in a trajectory. We start by extracting the interaction state-action segment of this trajectory. The interaction segment is converted to the embedding  $Z^A$  which guides the learning of  $z_i^{o,A}$  as explained before. The current state  $s_t$ , observation  $o_t$ , and observation embedding  $z_i^{o,A}$  determine what is the proper action  $a_t$  to take. The generation of  $z_i^{o,A}$  depends on an affordance-cue (aff-cue<sub>t</sub>) that is extracted from  $o_t$ . Note that the embedding  $z_i^{o,A}$  can essentially be viewed as a non-visual part of affordance-cue. However, in this paper, we focus on the visual part and we refer to the affordance-cue as an attention map.

## B. Siamese Encoder and Coupled Triplet Loss

Based on the graphical model, our design of Siamese encoder is depicted in the yellow region of Fig. 4. The Siamese encoder includes a LSTM layer to encode an interaction segment to generate  $Z_\tau^A$  ( $Z^A$  for a trajectory  $\tau$ ). It also includes a trajectory encoder that encodes all information (image, state, and previous action) per step to generate  $z_{i,\tau}^{o,A}$  ( $z_i^{o,A}$  at step  $t$  in a trajectory  $\tau$ ). Fig. 4 depicts the architecture at step  $t$ , but in the training phase, we feed into the Siamese encoder a trajectory of  $T$  steps and would obtain a sequence of observation embeddings  $\{z_{i,\tau}^{o,A}\}_{i=1}^T$ . We also replicate a Siamese encoder into three copies and they share the same weights at any time during training. The three copies take three trajectories as inputs during training: an anchor, positive, and negative trajectory. The anchor and positive trajectories are sampled from the same affordance category of data, while the anchor and negative trajectories come from two different categories. Given a trajectory, each copy of Siamese encoder generates  $Z_\tau^A$  and  $\{z_{i,\tau}^{o,A}\}_{i=1}^T$  for different possible  $\tau$  as explained before. Based on Equ. 1, 2 and 3, we propose the coupled triplet loss (Equ. 4) to couple the learning of the two types of embeddings together.

$$L(\mathcal{A}, \mathcal{P}, \mathcal{N}) = \sum_{i=1}^N \{ [\|Z_{\mathcal{A}_i}^A - Z_{\mathcal{P}_i}^A\|_2^2 - \|Z_{\mathcal{A}_i}^A - Z_{\mathcal{N}_i}^A\|_2^2 + M]_+ + \sum_{t=1}^T [ [\|Z_{\mathcal{A}_i}^{o,A} - z_{t,\mathcal{P}_i}^{o,A}\|_2^2 - \|Z_{\mathcal{A}_i}^{o,A} - z_{t,\mathcal{N}_i}^{o,A}\|_2^2 + M]_+ ] \} \quad (4)$$

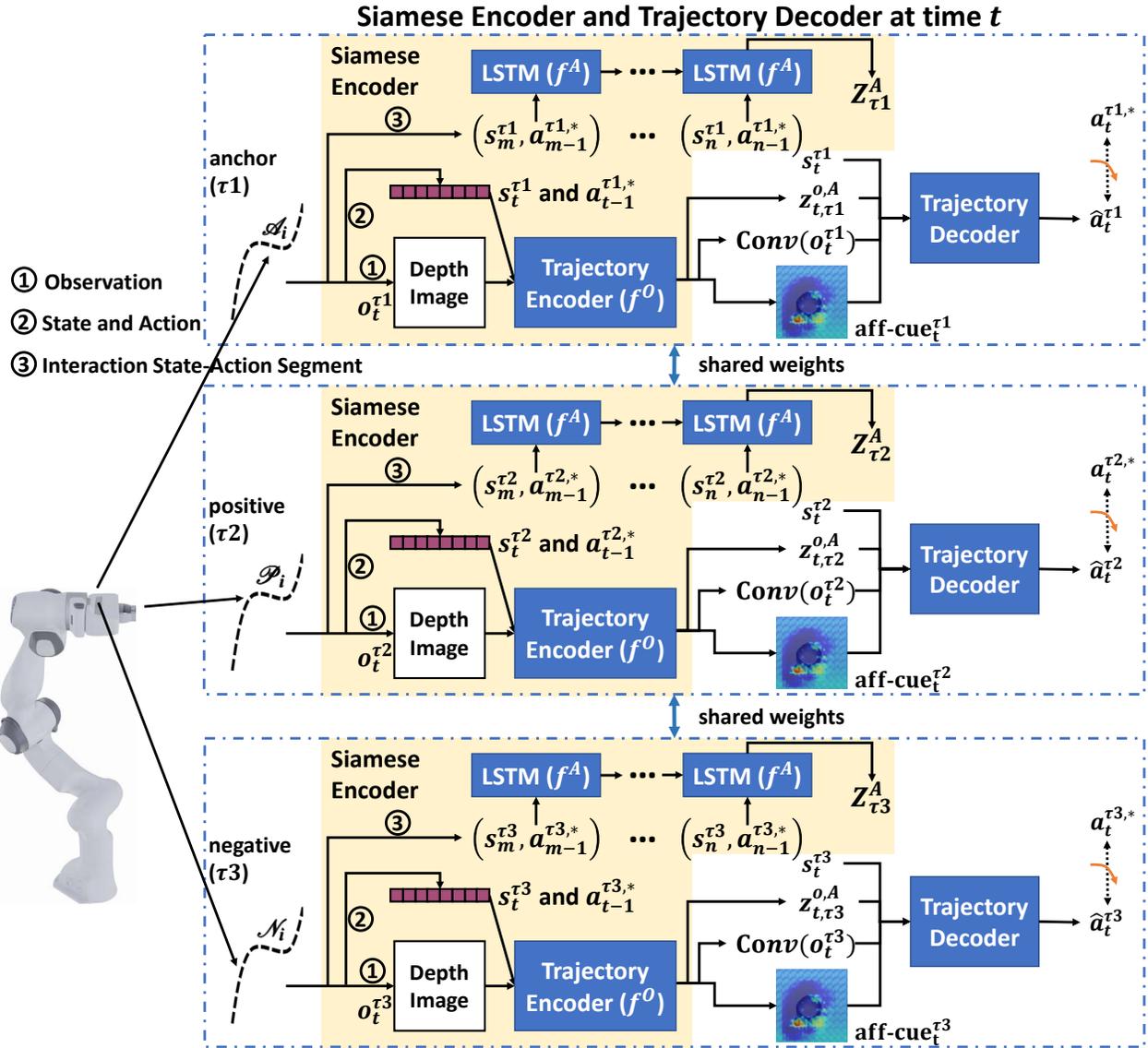


Fig. 4. The overall contrastive learning architecture of Siamese encoder and trajectory decoder. The dashed edge with an orange arrow means that the distance between its connected features needs to be minimized.

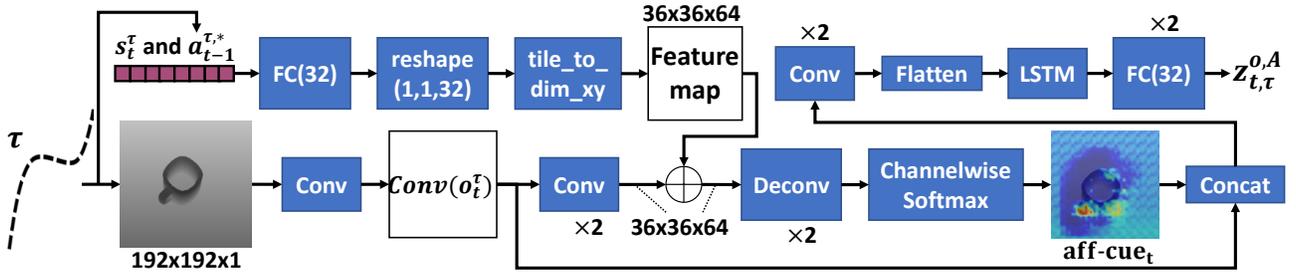


Fig. 5. The detailed architecture of the trajectory encoder module. “ $\times \#$ ” denotes that a neural network component needs to be replicated  $\#$  times.

where  $\mathcal{A}$ ,  $\mathcal{P}$ , and  $\mathcal{N}$  are the anchor, positive, and negative sets of demonstration trajectories;  $\mathcal{A}_i$  denotes the  $i$ -th trajectory in  $\mathcal{A}$  (likewise for  $\mathcal{P}_i$  and  $\mathcal{N}_i$ );  $T$  denotes the length of an arbitrary trajectory;  $f^A(\cdot)$  and  $f^O(\cdot)$  are explained under the Equ. 3, which generates an affordance embedding  $Z^A$  and an observation embedding  $z_t^{o,A}$  respectively.  $Z_\tau^A$  ( $\tau$  could

be either  $\mathcal{A}_i$ ,  $\mathcal{P}_i$ , or  $\mathcal{N}_i$ ) denotes an affordance embedding  $Z^A$  for a trajectory  $\tau$ ; Likewise,  $z_{i,\tau}^{o,A}$  denotes an observation embedding  $z_t^{o,A}$  for a trajectory  $\tau$ ; the Sec. IV-C explains  $f^O(\cdot)$  in more details.

The coupled triplet loss can be decomposed into two contrastive learning objectives in the two square brackets  $\square_+$

that are summed together. The first objective contrastively learns affordance embedding  $Z^A$  from all extracted interaction segments.  $Z^A$  could also be learned faster due to the low-dimensionality of interaction segment data. The second objective uses  $Z^A$  to guide the learning of the observation embedding  $z_t^{o,A}$ . Note that  $Z_{\mathcal{O}_i^A}^{o,A}$  is adjusted by both of the affordance embeddings  $z_{t,\mathcal{P}_i}^A$  and  $z_{t,\mathcal{A}_i}^A$ . This way of formulating the coupled triplet loss provides a strong training signal for the observation encoder that generates  $z_t^{o,A}$ . In this sense, the learning of the embeddings  $Z^A$  and  $z_t^{o,A}$  are coupled together.

### C. The Details of Siamese Encoder

The architecture of a Siamese encoder at an arbitrary time step is illustrated in the yellow region of Fig. 4 and with more details in Fig. 5. The input can be either an anchor, positive, or negative trajectory. The definition of a trajectory is explained in Sec. III.

We start by explaining the detailed formulation of the encoding function  $f^A()$ . At the initial step of an input trajectory, we append a dummy action vector of all zeros to provide a dummy previous action for the first step. We extract the interaction segment between the step  $m$  and  $n$  and feed them into an LSTM network ( $f^A()$ ) to generate the affordance encoding  $Z^A$ . In our work, the values of  $m$  and  $n$  are provided by a human expert. But in reality,  $m$  and  $n$  can be determined by tracking when the relative pose of the target object to gripper starts and stops changing.

We now explain the formulation of the encoding function  $f^O()$ . To process the observation encoding  $Z_t^{o,A}$  per step  $t$  of a whole trajectory, we feed it step-by-step into the trajectory encoder ( $f^O()$ ) module of Siamese encoder. As depicted in Fig. 5, the trajectory encoder module is a two-branch architecture. The bottom branch is used to process visual features and the top branch is used to process low-dimensional features like state and action vectors. The top branch concatenates state and action together and feeds them to a fully-connected layer. After that, it tiles (repeats) the fully-connected layer feature along the x and y dimensions (e.g.  $36 \times 36$ ) of the visual feature map generated after the third convolutional layer at the bottom branch. This way, the feature map representation of low-dimension inputs is merged with the convolutional features of visual inputs by element-wise addition. Then two deconvolution layers are used to generate a two-channel feature map. After applying the feature map with a channel-wise Softmax, its two channels represent graspable and non-graspable probabilities per pixel respectively. We then extract the first channel of this feature map as an attention map (affordance-cue),  $\text{aff-cue}_t$ . Note that in our work, such attention maps are essentially latent attention maps because their spatial dimension matches that of the first convolution layer ( $\text{Conv}(o_t^\tau)$ ) output, rather than the dimension of a raw input image. We then concatenate  $\text{aff-cue}_t$  with  $\text{Conv}(o_t^\tau)$ . After this, the two convolutional layers, one LSTM network, and two fully connected layers are used to generate  $Z_t^{o,A}$ . In our work we set the encoding size of  $Z_t^{o,A}$  to 32.

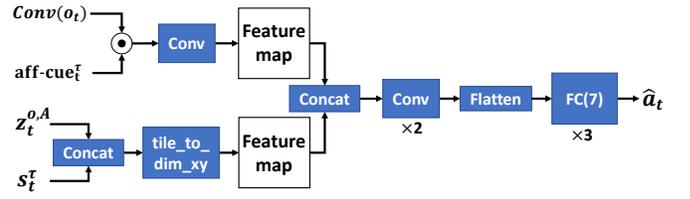


Fig. 6. The detailed architecture of the trajectory decoder. “ $\times\#$ ” denotes that a neural network component needs to be replicated  $\#$  times.

### D. Policy Network as Trajectory Decoder

The design of our trajectory decoder is based on a convolutional policy network, as illustrated in Fig. 6. The inputs include current state  $s_t$ , the latent convolution feature  $\text{Conv}(o_t)$ , the affordance cue ( $\text{aff-cue}_t$ ), and the contrastive embedding  $z_t^{o,A}$ . Since we treat a discovered affordance cue as an attention feature, we multiply  $\text{aff-cue}_t$  with each channel of the convolution feature  $\text{Conv}(o_t)$ .

We concatenate  $s_t$  and  $z_t^{o,A}$  together and obtain a new 1-D feature. We then tile (repeat) it across x and y dimensions of the visual feature map generated by the first convolutional layer at top branch. This way, we could concatenate the visual feature  $\text{Conv}(o_t)$ , state feature  $s_t$ , and contrastive embedding  $z_t^{o,A}$  together along the channel dimension. This new feature is then fed into two convolution layers and three fully connected layers to obtain a predicted action  $\hat{a}_t$ .

The loss function for behavior decoding is based on a behavior cloning loss in Equ. 5. The overall loss function for training the entire Siamese encoder and trajectory decoder is a sum of the coupled triplet loss (Equ. 4) and behavior cloning loss (Equ. 5), which enables a simultaneous learning of affordance knowledge and affordance-aware grasping from expert demonstrations.

$$\text{loss}_{bc} = \sum_{i=1}^N \left\{ \sum_{t=1}^T [L_1(a_t^*, \hat{a}_t) + L_2(a_t^*, \hat{a}_t)] \right\}_{\tau_i} \quad (5)$$

where  $L_1$  and  $L_2$  denotes L1-norm and L2-norm respectively;  $T$  denotes the length of a trajectory;  $|\tau_i$  denotes that the ground-truth action  $a_t^*$  is from the  $i$ -th trajectory  $\tau_i$  in dataset.

### E. Testing a Trained Model

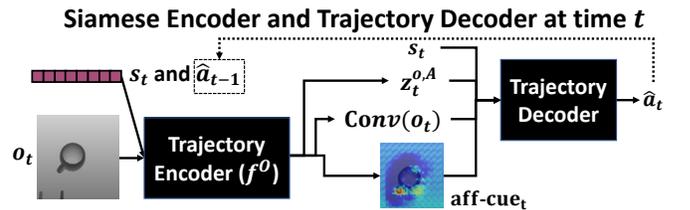


Fig. 7. The overall architecture of Siamese encoder and trajectory decoder in testing phase. The rectangular filled with black color means that the model weights are fixed. The dashed arrow and box represent that the prediction of action  $\hat{a}_t$  at each step would be fed into the Trajectory Encoder at the next step.

When we test a trained Siamese encoder and trajectory decoder (Fig. 7), we fix their neural network weights. Since

we do not start with an entire trajectory, there is no interaction segment that provides an affordance embedding. Instead, the trained weights in Trajectory Encoder already obtained affordance knowledge due to the coupled triplet loss (Equ. 4). Once the trajectory decoder outputs a predicted action  $\hat{a}_t$  at step  $t$ , the  $\hat{a}_t$ , rather than a ground-truth action as in Equ. 3, would be used for the computation at step  $t + 1$ .

## V. EVALUATION

We design our experiments in order to answer the following questions: 1) How helpful is the coupled triplet loss? 2) How helpful is the contrastive learning framework with a deep Siamese network? 3) How good is our model in comparison with a state-of-the-art baseline?

To answer 1, we compare our full model with a version that uses a normal triplet loss; to answer 2, we compare our full model with a version that does not have Siamese network and is trained totally based on behavior cloning losses; to answer 3, we compare our full model with a recently published work [20] as a baseline. Since we essentially solve an affordance-aware imitation learning problem for robotic grasping tasks, our evaluation metric involves grasping success rates which has been widely used for evaluating the learning of grasping/manipulation tasks (e.g. [12], [25]). We also show and analyze predicted affordance cues in a video from our project website<sup>1</sup> and partially in Fig. 1.

### A. Evaluation Domain and Data Collection

In our evaluation domain as shown in the leftmost region of Fig. 1, we have a mug that is put in front of a Franka Panda Arm in the PyBullet simulator [5]. We use 24 mugs in our experiments. These mugs have different affordance characteristics and each belongs to one or more of the three affordance categories. The task for the robot is to pick up the mug and lift it up for 5 centimeters. To do this, the robot needs to intelligently infer the best way of grasping the mug, e.g., by its body, the left and right sides of its handle, or the front and back sides of its handle. This domain reflects a common service that humans may need from robots in our everyday lives. This domain is also sufficient for evaluating our algorithm due to the fact that naturally there are a large number of mugs that have different structures and geometric characteristics. It could even be necessary to consider totally different ways of grasping them from different locations.

When we collect data, we randomly select a mug model and put it in a predefined position. We fix its initial pose across all of our experiments to guarantee the existence of an affordable grasp that can be categorized into either of our three affordance categories. All demonstration trajectories are 8 steps long. We use PyBullet’s function of reading users’ debugging commands to interactively move the gripper to a good target pose with a mouse. Once the gripper is moved to an ideal target pose, the robot then executes with that target pose, and relevant information like observation, action, and state are recorded. We collect 27 trajectories for our training dataset. After every 10 training epochs, we test

a model by randomly sampling 20 mugs and perform 20 grasps, respectively, and then record the success rate.

### B. Experimental Results and Analysis

The quantitative performance is measured by grasping success rates and the qualitative performance is evaluated by showing predicted affordance cues in our video as mentioned before. The grasping results of our model, the baseline, and ablation study models are reported in Table I. The success rate values are the highest testing success rates that a model achieves across all learning epochs. Table I provides experiment results of grasping success rates for four models: 1. our Siamese encoder with coupled triplet loss; 2. our Siamese network without coupled triplet loss (we apply normal triplet loss on observation embeddings); 3. our model that is not trained in a contrastive learning framework, and 4. a baseline behavior cloning work [20].

Model	Success Rate
1. Ours (Full Model): Siamese + Coupled Triplet Loss	65%
2. Ours (Ablation): Siamese + Normal Triplet Loss	35%
3. Ours (Ablation): Without Contrastive Learning	45%
4. Baseline [20]	25%

TABLE I

The results clearly show the advantage of using our coupled triplet loss with a Siamese neural network for the learning of affordance cues and grasping from demonstrations. An interesting finding is that the model 3 still performs better than model 2. This suggests that merely doing contrastive learning at observation level via the normal triplet loss could misguide the policy learning. Instead, the coupled triplet loss in this work contrastively learn affordance embeddings to guide the learning of observation encoding.

## VI. CONCLUSIONS AND FUTURE WORKS

We present an imitation learning algorithm that seeks training guidance not only from teachers’ actions but from simultaneously discovering teachers’ hidden affordance bias as well. We propose a contrastive learning framework with a Siamese encoder for affordance discovery and a trajectory decoder for policy learning. We represent affordance cues as visual attention. We further propose the coupled triplet loss to encourage the learned discriminative features to be more affordance-relevant. To the best of our knowledge, we initialize the direction of bridging the gap between affordance discovery and policy engineering/learning by achieving the two objectives together via an end-to-end deep neural network. Our work inherits the benefits of the class of works that learns affordances from demonstrations: there is no need of collecting ground-truth affordance labels for each image or pixel. However, such works focus on affordance learning and predicted affordances is still used by external motion planners. Our evaluation shows that our algorithm achieves the highest grasping success rate and predicts meaningful affordance cues. We believe our learning framework has a

larger potential for solving more complex tasks, which could be pursued in the future extensions of this work:

- 1) The tasks that involve multiple levels of interactions, instead of only the interaction between gripper and mug. Most tool-using tasks would require multiple levels of interactions;
- 2) The prediction of affordance-cue could be conditioned on different high-level tasks such that the attention map would be different even on the same object but w.r.t different tasks. In our work, we have a fixed high-level task of picking-to-pour-water, but if there is another task like pick-and-relocate, the discovered affordance cues might be different. This can be achieved by integrating our work into a hierarchical imitation learning framework.
- 3) It might also be interesting to explore how the coupled triplet loss could also be used to address other robot learning problems other than affordance-aware policy imitation.

#### REFERENCES

- [1] P. Abolghasemi, A. Mazaheri, M. Shah, and L. Boloni, "Pay attention! - robustifying a deep visuomotor policy through task-focused visual attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [2] S. J. Anderson, N. Yamagishi, and V. Karavia, "Attentional processes link perception and action," *Proceedings of the Royal Society of London. Series B: Biological Sciences*, vol. 269, no. 1497, pp. 1225–1232, 2002.
- [3] S. Banerjee, A. Hati, S. Chaudhuri, and R. Velmurugan, "Cosegnet: Image co-segmentation using a conditional siamese convolutional network," in *IJCAI*, 2019, pp. 673–679.
- [4] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah, "Signature verification using a siamese time delay neural network," in *Advances in neural information processing systems*, 1994, pp. 737–744.
- [5] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," 2016.
- [6] C. de Granville, J. Southerland, and A. H. Fagg, "Learning grasp affordances through human demonstration," in *Proceedings of the International Conference on Development and Learning (ICDL'06)*, 2006.
- [7] R. Detry, C. H. Ek, M. Madry, and D. Kragic, "Learning a dictionary of prototypical grasp-predicting parts from grasping experience," in *2013 IEEE International Conference on Robotics and Automation*. IEEE, 2013, pp. 601–608.
- [8] K. Fukunaga and P. M. Narendra, "A branch and bound algorithm for computing k-nearest neighbors," *IEEE transactions on computers*, vol. 100, no. 7, pp. 750–753, 1975.
- [9] J. J. Gibson, *The ecological approach to visual perception: classic edition*. Psychology Press, 2014.
- [10] K. Hsiao and T. Lozano-Perez, "Imitation learning of whole-body grasps," in *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2006, pp. 5657–5662.
- [11] T.-W. Huang, Y.-A. Wei, H.-T. Chen, and J. Liu, "Object discovery in depth images," in *2016 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*. IEEE, 2016, pp. 1–5.
- [12] K. Kleeberger, R. Bormann, W. Kraus, and M. F. Huber, "A survey on learning-based robotic grasping," *Current Robotics Reports*, pp. 1–11, 2020.
- [13] G. Koch, R. Zemel, and R. Salakhutdinov, "Siamese neural networks for one-shot image recognition," in *ICML deep learning workshop*, vol. 2. Lille, 2015.
- [14] K. Kostov and A. Janyan, "The role of attention in the affordance effect: can we afford to ignore it?" *Cognitive processing*, vol. 13, no. 1, pp. 215–218, 2012.
- [15] S. Li, S. Song, G. Huang, Z. Ding, and C. Wu, "Domain invariant and class discriminative feature learning for visual domain adaptation," *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4260–4273, 2018.
- [16] M. Lopes, F. S. Melo, and L. Montesano, "Affordance-based imitation learning in robots," in *2007 IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 2007, pp. 1015–1021.
- [17] X. Lu, W. Wang, C. Ma, J. Shen, L. Shao, and F. Porikli, "See more, know more: Unsupervised video object segmentation with co-attention siamese networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 3623–3632.
- [18] P. Mandikal and K. Grauman, "Dexterous robotic grasping with object-centric visual affordances," 09 2020.
- [19] S. May, M. Klodt, E. Rome, and R. Breithaupt, "Gpu-accelerated affordance cueing based on visual attention," in *2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2007, pp. 3385–3390.
- [20] J. Morton and M. J. Kochenderfer, "Simultaneous policy learning and latent state inference for imitating driver behavior," in *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2017, pp. 1–6.
- [21] P. Mukherjee, B. Lall, and S. Lattupally, "Object cosegmentation using deep siamese network," *arXiv preprint arXiv:1803.02555*, 2018.
- [22] E. Nepovimnykh, T. Eerola, and H. Kalviainen, "Siamese network based pelage pattern matching for ringed seal re-identification," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) Workshops*, March 2020.
- [23] K. Qian, X. Jing, Y. Duan, B. Zhou, F. Fang, J. Xia, and X. Ma, "Grasp pose detection with affordance-based task constraint learning in single-view point clouds," *Journal of Intelligent & Robotic Systems*, vol. 100, pp. 145–163, 2020.
- [24] K. Ramachandruni, M. Babu V., A. Majumder, S. Dutta, and S. Kumar, "Attentive task-net: Self supervised task-attention network for imitation learning using video demonstration," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 4760–4766.
- [25] A. Z. Ren, S. Veer, and A. Majumdar, "Generalization guarantees for multi-modal imitation learning," *arXiv preprint arXiv:2008.01913*, 2020.
- [26] L. Riggio, C. Iani, E. Gherri, F. Benatti, S. Rubichi, and R. Nicoletti, "The role of attention in the occurrence of the affordance effect," *Acta psychologica*, vol. 127, no. 2, pp. 449–458, 2008.
- [27] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 815–823.
- [28] L. Song, D. Gong, Z. Li, C. Liu, and W. Liu, "Occlusion robust face recognition based on mask learning with pairwise differential siamese network," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 773–782.
- [29] S. Srivastava, G. Sharma, and B. Lall, "Large scale novel object discovery in 3d," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018, pp. 179–188.
- [30] J. Steil, F. Röthling, R. Haschke, and H. Ritter, "Situated robot learning for multi-modal instruction and imitation of grasping," *Robotics and Autonomous Systems*, vol. 47, no. 2, pp. 129–141, 2004, robot Learning from Demonstration. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0921889004000430>
- [31] J. D. Sweeney and R. Grupen, "A model of shared grasp affordances from demonstration," in *2007 7th IEEE-RAS International Conference on Humanoid Robots*. IEEE, 2007, pp. 27–35.
- [32] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *European conference on computer vision*. Springer, 2016, pp. 499–515.
- [33] H. Yan, J. Lu, and X. Zhou, "Prototype-based discriminative feature learning for kinship verification," *IEEE Transactions on cybernetics*, vol. 45, no. 11, pp. 2535–2545, 2014.
- [34] L. Yen-Chen, A. Zeng, S. Song, P. Isola, and T.-Y. Lin, "Learning to see before learning to act: Visual pre-training for manipulation," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2020. [Online]. Available: <https://yenchenlin.me/vision2action/>
- [35] A. Zeng, S. Song, K.-T. Yu, E. Donlon, F. R. Hogan, M. Bauza, D. Ma, O. Taylor, M. Liu, E. Romo, et al., "Robotic pick-and-place of novel objects in clutter with multi-affordance grasping and cross-domain image matching," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 3750–3757.
- [36] X. Zhang, F. Pala, and B. Bhanu, "Attributes co-occurrence pattern mining for video-based person re-identification," in *2017 14th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*. IEEE, 2017, pp. 1–6.