



AdaAfford: Learning to Adapt Manipulation Affordance for 3D Articulated Objects via Few-Shot Interactions

Yian Wang^{1,2}, Ruihai Wu^{1,2}, Kaichun Mo³, Jiaqi Ke^{1,2}, Qingnan Fan⁴,
Leonidas J. Guibas³, and Hao Dong^{1,2,5}(✉)

¹ CFCS, CS Department, PKU, Beijing, China
{yianwang,wuruihai,kjq001220,hao.dong}@pku.edu.cn
² AIT, PKU, Beijing, China
³ Stanford University, Stanford, USA
{kaichun,guibas}@cs.stanford.edu
⁴ Tencent AI Lab, Bellevue, USA
⁵ Peng Cheng Lab, Shenzhen, China
<https://hyperplane-lab.github.io/AdaAfford>

Abstract. Perceiving and interacting with 3D articulated objects, such as cabinets, doors, and faucets, pose particular challenges for future home-assistant robots performing daily tasks in human environments.

Besides parsing the articulated parts and joint parameters, researchers recently advocate learning manipulation affordance over the input shape geometry which is more task-aware and geometrically fine-grained.

However, taking only passive observations as inputs, these methods ignore many hidden but important kinematic constraints (*e.g.*, joint location and limits) and dynamic factors (*e.g.*, joint friction and restitution), therefore losing significant accuracy for test cases with such uncertainties. In this paper, we propose a novel framework, named AdaAfford, that learns to perform very few test-time interactions for quickly adapting the affordance priors to more accurate instance-specific posteriors. We conduct large-scale experiments using the PartNet-Mobility dataset and prove that our system performs better than baselines. We will release our code and data upon paper acceptance.

1 Introduction

For future home-assistant robots to aid humans in accomplishing diverse everyday tasks, we must equip them with strong capabilities perceiving and interacting with diverse 3D objects in human environments. Articulated objects, such

Y. Wang, R. Wu and K. Mo—Equal contribution.

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-19818-2_6.

as cabinets, doors, and faucets, are particularly interesting kinds of 3D shapes in our daily lives since agents can interact with them and trigger functionally important state changes of the objects (*e.g.*, push closed the drawer of the cabinet, rotate the handle and pull open the door, turn on/off the water from the faucet by rotating the switch). However, because robots need to understand more semantically complicated part semantics and manipulate articulated parts with higher degree-of-freedom than rigid objects, it remains a very important yet challenging task to perceive and interact with 3D articulated objects.

Many previous works have investigated the problem of perceiving and interacting with 3D articulated objects. Researchers have been pushing the state-of-the-arts on segmenting articulated parts [32, 42], tracking them [30, 35], and estimating joint parameters [34, 40], enabling robotic systems [2, 23, 33] to successfully perform sophisticated planning and control over 3D articulated objects.

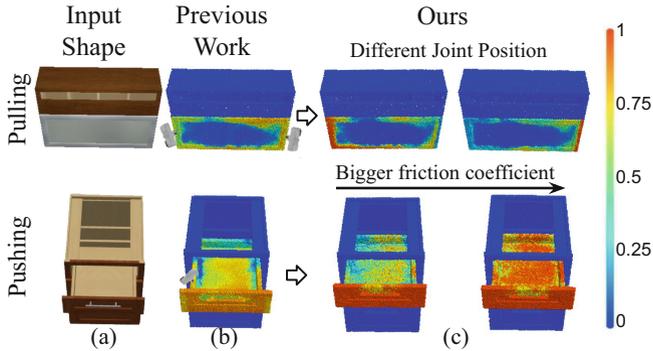


Fig. 1. For robotic manipulation over 3D articulated objects (a), past works [18, 36] have demonstrated the usefulness of per-point manipulation affordance (b). However, only observing static visual inputs passively, these systems suffer from intrinsic ambiguities over kinematic constraints. Our *AdaAfford* framework reduces such uncertainties via interactions and quickly adapts instance-specific affordance posteriors (c).

More recently, beyond recognizing the articulated parts and joints, researchers have been proposing learning more task-aware and geometrically fine-grained manipulation affordance over input 3D geometry. Where2Act [18], the most related to our work, learns densely labeled manipulation affordance heatmaps over 3D input partial scans of articulated objects, as illustrated in Fig. 1(b), by performing self-supervised trial-and-error interaction in a physical simulator. There are also many other works leveraging similar dense affordance predictions over 3D scenes [21] and rigid objects [17]. Such densely labeled affordance predictions over 3D data provide more geometrically fine-grained actionable information and can be learned task-specifically given different manipulation actions, showing promises in bridging the perception-interaction gaps for robotic manipulation over large-scale 3D data across different tasks.

However, taking only a single-frame observation of the 3D shape as input (*e.g.*, a single 2D image, a single partial 3D scan), these methods systematically fail to capture many hidden but important kinematic or dynamic factors and therefore predict inaccurate affordance heatmaps, similar to Fig. 1(b), by averaging out such uncertainties. For example, given a fully closed cabinet door with no obvious handle as shown in Fig. 1 (top-row), it is uncertain if the door axis is on the left or right side, which significantly affects the manipulation affordance predictions. Other kinematic uncertainties include joint limits (*e.g.*, push inward or pull outward for a door) and types (*e.g.*, slide or rotate to open a door). Besides, various dynamic or physical parameters (*e.g.*, part mass, joint friction) are also unobservable from single-frame inputs but largely affect manipulation affordance. For example, with increasing friction coefficient for a cabinet drawer (Fig. 1, bottom-row), robots would be able to push the inner board.

In this paper, we propose a novel framework *AdaAfford* learning perform very few test-time interactions to reduce such kinematic or dynamic uncertainties and fastly adapts the affordance prior predictions to instance-specific posteriors given a novel test shape. Our system learns a *data-efficient strategy* that sequentially samples very few uncertain or interesting locations to interact, as the interacting grippers illustrated in Fig. 1(b), according to the current affordance predictions and past interaction trials (we begin with the affordance prior predictions of Where2Act [18] and zero interaction history). The interaction outcomes, each of which includes the interaction location, direction, and the resulting part motion, are then observed and incorporated to produce posterior affordance predictions, as illustrated in Fig. 1(c), by a proposed *fast-adaptation mechanism*.

We set up a benchmark for experiments and evaluations using the large-scale PartNet-Mobility dataset [20] and the SAPIEN physical simulator [37]. We use in total 972 shapes from 15 object categories and conduct experiments for several action types, and randomly sample the kinematic and dynamic parameters for the 3D articulated objects in simulation. Experiments show our method can successfully and efficiently adapt manipulation affordance to novel test shapes with as few as one to four interactions. Quantitative evaluation further proves the effectiveness of our proposed approach.

In summary, our main contributions are the following. 1) we point out and investigate an important limitation of the methods that learn densely labeled visual manipulation affordance – the unawareness of hidden yet important kinematic and dynamic uncertainties; 2) we propose a novel framework *AdaAfford* that learns to perform very few test-time interactions to reduce uncertainties and quickly adapt to predicting an instance-specific affordance posterior; 3) we set up a large-scale benchmark, built upon PartNet-Mobility [20] and SAPIEN [37], for experiments and evaluations, and results demonstrated the effectiveness and efficiency of the proposed approach.

2 Related Work

Visual Affordance on 3D Shapes. Affordance [9] suggests possible ways for agents to interact with objects. Many past works have investigated

learning grasp [11, 13, 15, 26, 29] and manipulation [17, 18, 21, 27, 36, 39] affordance for robot-object interaction, while there are also many works studying affordance for hand-object [3, 4, 12, 17, 41], object-object [19, 31, 46], and human-scene [8, 16, 21, 25] interaction scenarios. Among these works, researchers have proposed different representations for visual affordance, including detection locations [15, 29], parts [17], keypoints [27], heatmaps [18, 21], etc. In this work, we mostly follow the settings in [18] for learning visual affordance heatmaps for manipulating 3D articulated objects. Different from previous works that infer possible agent-object visual affordance heatmaps passively from static visual observations, our framework leverages active interactions to efficiently query uncertain kinematic or dynamic factors for learning more accurate instance-adaptive visual affordance.

Fast Adaption via Few-Shot Interactions. Researchers have explored various approaches [5, 7, 28, 44, 45] for fast adaption via few-shot interactions. Many past works have also designed interactive perception methods to figure out object mass [14], dynamic parameters [1, 6, 10, 38], or parameters for known models [43]. Different from these studies proposing general algorithms for policy adaptation or figuring out explicit system parameters for rigid objects, we focus on designing a working solution for our specific task of learning visual affordance heatmaps for manipulating 3D articulated objects with special designs on predicting geometry-grounded interaction proposals and interaction-adaptive affordance predictions.

3 Problem Formulation

Given as input a single-frame 3D partial point cloud observation of an articulated object $O \in \mathbb{R}^{N \times 3}$ (*e.g.*, lifted from a depth scanner with known camera intrinsics), the Where2Act framework [18] directly outputs a per-point manipulation affordance heatmap $A \in [0, 1]^N$, where higher scores indicate bigger chances for being interacted with to accomplish a given short-term manipulation task (*e.g.*, pushing, pulling). Additionally, a diverse set of gripper orientations $\{R_1^p, R_2^p, \dots, R_i^p \in SO(3)\}$ is proposed at each point $p \in O$ suggesting possible ways for robot agents to interact with, each of which also associated with a success likelihood $s_i^p \in [0, 1]$. No interaction is allowed at test time in Where2Act and a fixed set of system dynamic parameters is used across all shapes.

We follow most of the Where2Act settings except that we randomly vary the system dynamics and allow test-time interactions over the 3D shape to reduce kinematic or dynamic uncertainties. Our AdaAfford system proposes a few interactions sequentially $\mathcal{I} = \{I_1, I_2, \dots\}$. Each interaction $I_i = (O_i, p_i, R_i, m_i)$ executes a task-specific hard-coded short-term trajectory defined in Where2Act, parametrized by the interaction point $p_i \in O_i$ and the gripper orientation $R_i \in SO(3)$, and observes a part motion m_i . Starting from the input shape observation $O_1 \leftarrow O$, every interaction I_i where $m_i \neq 0$ changes the part state and thus produces a new shape point cloud input for the next interaction $O_{i+1} \neq O_i$. Leveraging the interaction observations \mathcal{I} , our system then adapts the per-point manipulation affordance A predicted by Where2Act to a posterior $A_{\mathcal{I}} \in [0, 1]^N$

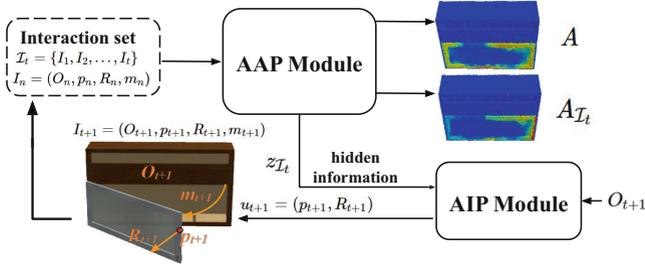


Fig. 2. Method Overview. Starting from the Where2Act [18] predicted affordance prior A , at each timestep $t = 1, 2, \dots$, we recursively leverage the *Adaptive Interaction Proposal* (AIP) module to propose a next-time interaction action u_{t+1} , observe the interaction outcome m_{t+1} , and feed through the *Adaptive Affordance Prediction* (AAP) module all past few-shot interactions \mathcal{I}_t together with the new one I_{t+1} for adapting to an affordance posterior prediction $A_{\mathcal{I}_{t+1}}$. The procedure iterates until the interaction budget is reached or the AIP module decides to stop.

that reduces uncertainties and provides more accurate instance-specific predictions. For each gripper orientation R_i , we also update the success likelihood score $s_{i,\mathcal{I}}^p \in [0, 1]$ considering the test-time interactions.

4 Method

Our proposed *AdaAfford* framework primarily consists of two modules – an *Adaptive Interaction Proposal* (AIP) module and an *Adaptive Affordance Prediction* (AAP) module. While the AIP module learns a greedy yet effective strategy for sequentially proposing few-shot test-time interactions $\mathcal{I} = \{I_1, I_2, \dots\}$ revealing hidden information, the AAP module is trained to adapt affordance predictions from Where2Act [18] prior A to a posterior $A_{\mathcal{I}}$ observing the sampled interactions \mathcal{I} . We iterate two modules recurrently at test time to produce a sequence of few-shot interactions \mathcal{I} leading to the final affordance posterior prediction $A_{\mathcal{I}}$. During training, we iteratively alternate the training for the two modules until a joint convergence. Below, we first introduce the test-time inference procedure for a brief overview. Next, we describe the input backbone encoders that are shared among all networks in our framework. Then, we describe the detailed architectures and system designs of the two modules. We conclude with the training losses and strategy.

Test-Time Overview. Figure 2 presents an overview of the method. We apply a recurrent structure at test time. Starting from the affordance prediction A , the AIP module proposes the first action for producing the interaction data I_1 . Then, at each timestep $t = 1, 2, \dots$, we feed the current set of interactions $\mathcal{I}_t = \{I_1, \dots, I_t\}$ as inputs to the AAP module and extract hidden information $z_{\mathcal{I}_t} \in \mathbb{R}^{128}$ that adapts the affordance map prediction to $A_{\mathcal{I}_t}$. The AIP module then takes $z_{\mathcal{I}_t}$ as input and proposes an action $u_{t+1} = (p_{t+1}, R_{t+1})$ composed of the

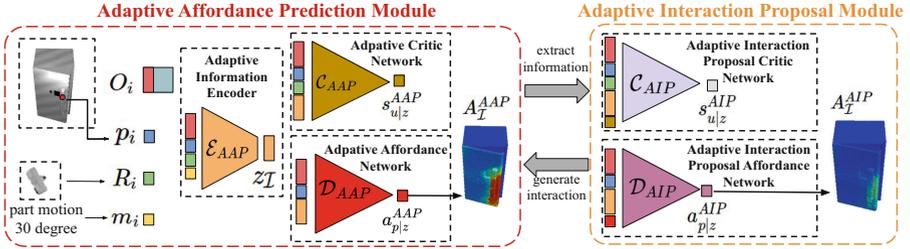


Fig. 3. Network Architecture. Left: the *Adaptive Affordance Prediction* (AAP) module takes as inputs the few-shot interactions \mathcal{I} and predicts the affordance posterior $A_{\mathcal{I}}$. Right: the *Adaptive Interaction Proposal* (AIP) module proposes a next-step interaction action $u_{t+1} = (p_{t+1}, R_{t+1})$ (denote the current timestep as t) given the feature $z_{\mathcal{I}}$ extracted from the current interaction observations \mathcal{I} .

interaction point p_{t+1} and the gripper orientation R_{t+1} for the next interaction. Performing this action in the environment, we obtain the next-step interaction data $I_{t+1} = (O_{t+1}, p_{t+1}, R_{t+1}, m_{t+1})$ and put it into the interaction set $\mathcal{I}_{t+1} \leftarrow \mathcal{I}_t \cup \{I_{t+1}\}$. We iterate until the interaction budget has been reached or our AIP module decides to stop. When the procedure stops at timestep T , we output the final affordance posterior $A_{\mathcal{I}} = A_{\mathcal{I}_T}$.

Input Encoders. This paragraph details how we encode inputs into features as all the encoder networks in the two modules take the same input entities (*e.g.*, the shape observation O , the interaction action u) and thus share the same architecture. We use the PointNet++ segmentation network [24] to encode the input shape point cloud $O \in \mathbb{R}^{N \times 3}$ into per-point feature maps $f_O \in \mathbb{R}^{N \times 128}$ and denote $f_{p|O} \in \mathbb{R}^{128}$ as the feature at any point $p \in O$. We use Multilayer Perceptron (MLP) networks to encode other vector inputs (*e.g.*, the interaction action u and the part motion m) into $f_a \in \mathbb{R}^{128}$. The networks in the following subsections will first encode the inputs into $f_{p|O}$ and f_a , and then concatenate them into $f_I \in \mathbb{R}^{256}$. The encoders do not share weights across different modules (Fig. 3).

4.1 Adaptive Affordance Prediction Module

The *Adaptive Affordance Prediction* (AAP) module takes as inputs few-shot interactions \mathcal{I} and predicts the affordance posterior $A_{\mathcal{I}}$. This module is composed of three subnetworks: 1) an *Adaptive Information Encoder* \mathcal{E}_{AAP} that extracts hidden information $z \in \mathbb{R}^{128}$ from a set of interactions \mathcal{I} ; 2) an *Adaptive Affordance Network* \mathcal{D}_{AAP} that predicts the posterior affordance heatmap $A_{\mathcal{I}}$ conditioned on the hidden information z ; and 3) an *Adaptive Critic Network* \mathcal{C}_{AAP} that predicts the AAP action score $s_{u|z}^{AAP} \in [0, 1]$ for an action u conditioned on the hidden information z . Here, an action is represented as $u = (p, R)$ including an interaction point $p \in O$ and a gripper orientation $R \in SO(3)$.

Adaptive Information Encoder. Given a set of interactions $\mathcal{I} = \{I_1, I_2, \dots\}$ as inputs, the *Adaptive Information Encoder* \mathcal{E}_{AAP} outputs a 128-dim hidden information representation $z_{\mathcal{I}}$ (z for brevity). It first encodes each interaction I_i using the input encoders mentioned before, and then uses an MLP network to encode the features into a 128-dim latent code z_{I_i} representing the hidden information extracted from I_i . As different interactions contain different amount of hidden information, we use another MLP Network to predict an attention score $w_{I_i} \in \mathbb{R}$ for each interaction. To get a summarized hidden information from a set of interactions, we simply compute a weighted average over all z_{I_i} 's according to the weights w_{I_i} 's and use the resulting feature as $z_{\mathcal{I}}$. Formally, we have $z_{\mathcal{I}} \leftarrow (\sum_i z_{I_i} \times w_{I_i}) / (\sum_i w_{I_i})$.

Adaptive Critic Network. Given the object partial point cloud observation O , an arbitrary interaction point $p \in O$, an arbitrary gripper orientation $R \in SO(3)$ and the latent code z , the *Adaptive Critic Network* \mathcal{C}_{AAP} predicts an AAP action score $s_{u|z}^{AAP} \in [0, 1]$ indicating the likelihood for the success of the interaction action u given the interaction information z . It first encodes the input $\{O, P, R\}$ using the input encoders as mentioned before and then employs an MLP network to predict AAP action score $s_{u|z}^{AAP}$, taking the concatenated features together with z as inputs. A higher AAP action score $s_{u|z}^{AAP}$ for action u indicates a higher chance for u to succeed in accomplishing the given manipulation task.

Adaptive Affordance Network. Given the input object partial point cloud O , an arbitrary point $p \in O$, and the latent code z , the *Adaptive Affordance Network* \mathcal{D}_{AAP} predicts an actionability score $a_{p|z}^{AAP} \in [0, 1]$ at point p . It first encodes the input $\{O, p\}$ using the aforementioned input encoders and then uses an MLP network that takes the concatenated features together with z as inputs and produces an actionability score $a_{p|z}^{AAP}$ as the output. A higher actionability score $a_{p|z}^{AAP}$ indicates a higher chance to successfully interact on point p .

4.2 Adaptive Interaction Proposal Module

Adaptive Interaction Proposal. (AIP) module proposes an action (denote the current timestep as t) $u_{t+1} = (p_{t+1}, R_{t+1})$ for the next step interaction, given the feature z extracted from the current interaction observations \mathcal{I} . This module contains two networks: 1) an *Adaptive Interaction Proposal Affordance Network* \mathcal{D}_{AIP} that predicts an AIP actionability score $a_{p|z}^{AIP} \in \mathbb{R}$ indicating how likely the next-action is worth interacting at point p , and 2) an *Adaptive Interaction Proposal Critic Network* \mathcal{C}_{AIP} predicting an AIP action score $s_{u|z}^{AIP} \in \mathbb{R}$ suggesting the gripper orientation to pick for the next interaction. We leverage the predictions of the two networks to propose the next action $u_{t+1} = (p_{t+1}, R_{t+1})$.

Adaptive Interaction Proposal Critic Network. Given the input object partial point cloud O , an arbitrary interaction point $p \in O$, an arbitrary gripper orientation $R \in SO(3)$, the latent code z , and the AAP action score $s_{u|z}^{AAP}$ produced by \mathcal{C}_{AAP} , the *AIP Critic Network* \mathcal{C}_{AIP} predicts the AIP action score

$s_{u|z}^{AIP} \in \mathbb{R}$ of u . It first encodes the inputs $\{O, p, R, s_{u|z}^{AAP}\}$ using the input encoders and then uses an MLP network that takes the concatenated features together with z as inputs and generates an AIP action score $s_{u|z}^{AIP}$ for the action u . A higher AIP action score suggests that the action u may query more unknown yet interesting hidden information and thus is worth exploring next.

Adaptive Interaction Proposal Affordance Network. Given the input partial shape observation O , an arbitrary interaction point $p \in O$, the latent code z , and the AAP actionability score $a_{p|z}^{AAP}$ at point p estimated by \mathcal{D}_{AAP} , the *AIP Affordance Network* \mathcal{D}_{AIP} predicts the AIP actionability score $a_{p|z}^{AIP} \in \mathbb{R}$ at point p . It first encodes the inputs $\{O, p, a_{p|z}^{AAP}\}$ using the aforementioned input encoders and then employs an MLP network to predict an AIP actionability score $a_{p|z}^{AIP}$, taking the concatenated features together with z as inputs. A higher AIP actionability score at p indicates more unknown yet helpful hidden information may be obtained by executing an interaction at p .

Next-Step Interaction Action Proposal. In order to propose an action $u_{t+1} = (p_{t+1}, R_{t+1})$ for the next interaction, given the hidden information z and the input shape partial point cloud O , we first obtain the AIP actionability heatmap $A_{p|z}^{AIP}$ for every point $p \in O$ predicted by the *AIP Affordance Network* \mathcal{D}_{AIP} and then select the point $p_{t+1} \leftarrow p_*$ with the highest AIP actionability score $a_{p_*|z}^{AIP}$. Then, we sample 100 random actions $\{u_1, u_2, \dots, u_{100}\}$ at p using the Where2Act’s pre-trained *Action Proposal Network*, use our *AIP critic network* \mathcal{C}_{AIP} to generate the AIP action scores $s_{u_i|z}^{AIP}$ for each action u_i , and then choose the action $u_{t+1} \leftarrow u_*$ with the highest AIP action score $s_{u_*|z}^{AIP}$.

Stopping Criterion for the Few-Shot Interactions. The AIP procedure for generating few-shot interactions stops when a preset budget is reached or the maximal AIP actionability score is below a certain threshold (e.g., 0.05).

4.3 Training and Losses

In brief, for AAP module, we use ground-truth motion m to supervise \mathcal{E}_{AAP} and \mathcal{C}_{AAP} , and utilize \mathcal{C}_{AAP} to supervise the training of \mathcal{D}_{AAP} . For AIP module, we use AAP module to supervise the training of \mathcal{C}_{AIP} and use it to supervise \mathcal{D}_{AIP} . Below, we describe the losses and the training strategy in detail.

AAP Action Scoring Loss. To supervise \mathcal{C}_{AAP} , we use a standard binary cross entropy loss, which measures the error between the prediction of \mathcal{C}_{AAP} and target part’s ground truth motion m of an interaction I . Specifically, given the hidden information z , a batch of interaction observations $\mathcal{I} = \{I_1, I_2, \dots, I_B\}$ where $I_i = \{O_i, u_i, m_i\}$, and the AAP action score prediction $s_{u_i|z}^{AAP}$ for each interaction I_i , the loss is defined as

$$\mathcal{L}_C^{AAP} = -\frac{1}{B} \sum_i r_i \log(s_{u_i|z}^{AAP}) + (1 - r_i) \log(1 - s_{u_i|z}^{AAP})$$

where $r_i = 1$ if $m_i > \tau$ (e.g., $\tau = 0.01$) or $r_i = 0$ rendering a binary discretization for each interaction outcome.

AAP Actionability Scoring Loss. To train \mathcal{D}_{AAP} , we apply an \mathcal{L}_1 loss to measure the difference from the predicted score $a_{p|z}^{AAP}$ to the ground truth. To estimate the ground truth actionability score for p , we randomly sample 100 actions at p according to pre-trained Where2Act *Action Proposal Network*, predict AAP action scores $s_{u|z}^{AAP}$'s of these actions u 's using \mathcal{C}_{AAP} , and take the average of the top-5 scores as the ground truth actionability score.

AIP Action Scoring Loss. To supervise \mathcal{C}_{AIP} , we use an \mathcal{L}_1 loss to measure the difference between our predicted AIP action score $s_{u|z}^{AIP}$ and the ground truth AIP action score $gt_{u|z}^{AIP}$. Given a set of interactions $\mathcal{I}_T = \{I_1, I_2, \dots\}$, to generate $gt_{u_i|z}^{AIP}$ for an interaction action u_i , we respectively encode two interaction subsets $\mathcal{I}_{i-1} = \{I_1, I_2, \dots, I_{i-1}\}$ and $\mathcal{I}_i = \{I_1, I_2, \dots, I_i\}$ into latent codes $z_{\mathcal{I}_i}$ and $z_{\mathcal{I}_{i-1}}$. Then feed $z_{\mathcal{I}_i}$ and $z_{\mathcal{I}_{i-1}}$ as the conditional inputs to \mathcal{C}_{AAP} separately and count the difference of $\mathcal{L}_{\mathcal{C}^{AAP}}$ as the ground truth of AIP action score $gt_{u_i|z_{\mathcal{I}_{i-1}}}^{AIP}$. More concretely, let the AAP action scoring loss conditioned on $z_{\mathcal{I}_i}$ and $z_{\mathcal{I}_{i-1}}$ respectively be $\mathcal{L}_{\mathcal{I}_i}$ and $\mathcal{L}_{\mathcal{I}_{i-1}}$. We define the ground truth AIP action score $gt_{u_i|z_{\mathcal{I}_{i-1}}}^{AIP} \leftarrow \mathcal{L}_{\mathcal{I}_{i-1}} - \mathcal{L}_{\mathcal{I}_i}$. The AIP action score is trained to regress an estimated positive influence of executing u on the AAP action score predictions, where an action giving more influence is preferred as it helps discover more hidden information useful to the task.

AIP Actionability Scoring Loss. To train \mathcal{D}_{AIP} , we use another \mathcal{L}_1 loss. For each $p \in O$, we sample 100 actions u_i 's using the pre-trained Where2Act *Action Proposal Network*, obtain the AIP action scores $s_{u_i|z}^{AIP}$'s of these actions u_i 's by \mathcal{C}_{AIP} , and use the average of the top-5 scores as the regression target.

Training Strategy. We iteratively train the AAP module and AIP module until a joint convergence since the update of the subnetworks in one module will affect the training of the subnetworks in the other module. More specifically, the update of \mathcal{C}_{AAP} in the AAP module will affect the ground-truth AIP action scores, while the update of \mathcal{C}_{AIP} and \mathcal{D}_{AIP} in the AIP module will change the proposed interactions used to generate z in the AAP module. Therefore, our final solution is to train the AAP and AIP modules iteratively.

5 Experiments

We perform experiments using the large-scale PartNet-Mobility dataset [20] and the SAPIEN simulator [37], and set up several baselines for comparisons. Results demonstrate the effectiveness and superiority of the proposed approach.

5.1 Data and Settings

Data. Following the settings of Where2Act [18], we conduct our experiments in the SAPIEN [37] simulator equipped with NVIDIA PhysX [22] simulation engine

and the large scale PartNet-Mobility [20] dataset. We use 972 articulated 3D objects covering 15 object categories, mostly following Where2Act, to carry out the experiments. The dataset is divided into 10 training and 5 testing categories. The shapes in the training categories are further divided into two disjoint sets of training and test shapes. See supplementary for detailed statistics.

Table 1. Quantitative Evaluations. We experiment with three different test-time interaction budgets (*i.e.*, 1, 2, or 4) where numbers are separated by slashes. We use “pushing all” and “pulling all” to denote the experiments over all object categories, while “pulling closed door” and “pushing faucet” refer to the experiments over a single category only. For the experiments over all categories, we report the performance over novel shapes from the training categories (marked with “train cat.”) and shapes from novel categories (marked with “test cat.”).

		F-score (%)	Sample-Succ (%)
Pushing all (train cat.)	Where2Act	56.44	20.85
	Where2Act-adaptation	64.16/65.42/64.99	20.77/22.72/26.82
	Ours-fps	64.32/69.58/70.99	26.22/27.30/30.65
	Ours-final	72.78/73.12/75.18	33.82/33.23/35.23
Pushing all (test cat.)	Where2Act	59.95	21.69
	Where2Act-adaptation	51.09/53.28/55.56	19.06/22.27/24.50
	Ours-fps	66.17/67.27/69.08	33.64/35.19/ 37.79
	Ours-final	77.58/77.63/78.42	34.97/36.75/37.40
Pulling all (train cat.)	Where2Act	31.19	1.92
	Where2Act-adaptation	37.22/38.48/39.13	1.11/2.15/1.62
	Ours-fps	39.88/42.74/43.55	2.78/5.56/4.44
	Ours-final	42.62/43.87/44.08	7.78/9.44/10.55
Pulling all (test cat.)	Where2Act	36.36	10.00
	Where2Act-adaptation	40.11/45.52/48.80	3.40/6.25/10.17
	Ours-fps	43.67/42.77/48.33	4.35/3.91/4.78
	Ours-final	49.51/50.00/51.33	5.21/7.39/10.45
Pulling closed door	Where2Act	48.44	4.38
	Where2Act-adaptation	50.21/55.75/56.81	6.60/7.18/6.83
	Ours-fps	59.79/63.43/69.13	8.88/11.33/12.10
	Ours-final	57.83/ 65.60/79.65	10.86/11.57/22.14
Pushing faucet	Where2Act	64.92	55.46
	Where2Act-adaptation	66.25/62.18/67.15	57.50/52.08/61.70
	Ours-fps	74.19/79.36/77.95	60.44/70.12/77.41
	Ours-final	77.42/83.06/83.83	65.90/81.66/82.14

Experiment Settings. Following Where2Act [18], we perform experiments over all object categories under different manipulation action types. We train one network for each downstream manipulation task over training shapes from the 10 training object categories and evaluate the performance over test shapes from the training categories and shapes from unseen test categories. Besides, to further demonstrate the effectiveness of our method, we conduct two additional experiments under challenging tasks with clear kinematic ambiguity, each of which is conducted over a single object category: 1) pulling closed doors of cabinets that cannot be easily distinguished which side to pull open; 2) pushing faucets with uncertainties which direction to rotate (clockwise or counter-clockwise). These experiments are particularly interesting yet challenging cases on which previous work Where2Act [18] fail drastically and we hope to test our framework.

Environment Settings. Following Where2Act, we abstract away the robot arm and only use a Franka Panda flying gripper as the robot actuator. The input shape point cloud is assumed to be cleanly segmented out. To generate the input partial point cloud scans, we mount an RGB-D camera with known intrinsic parameters 5-unit-length away pointing to the center of the target object.

To simulate manipulating shapes with uncertain dynamics, we randomly change the following three physical parameters in SAPIEN: 1) the friction of the target part joint, 2) the mass of the target part, and 3) the friction coefficient of the target part surface. For the “pulling closed door” task, we manually select the cabinets whose doors have no clear handle geometry in the PartNet-Mobility dataset [37], and set the poses of those doors to be closed. The gripper cannot tell which side to pull open the door because it is impossible to tell whether the axis position is on the left or right of the door from passive visual observations. For the “pushing faucet” task, we randomly set the rotating direction of the faucet switch to be in one of the following three modes: only clockwise, only counter clockwise, or both ways.

5.2 Baselines and Evaluation Metrics

We set up several baseline and employ two metrics for quantitative comparisons.

Baselines and Ablation Study. We compare our framework with several baselines (see supplementary for more detailed descriptions for the baseline designs):

- **Where2Act:** the original method proposed in [18] where only the pure visual information is used for predicting the visual actionable information and no interaction data is used at all during test time;
- **Where2Act-adaptation:** the Where2Act method augmented with a heuristic based adaptation mechanism to replace the AAP module where given the interaction observations we locally adjust the predictions for similar points;
- **Ours-fps:** a variant of our proposed method that we use FPS to sample over the predicted affordance for interactions instead of the AIP proposals.

We compare to **Where2Act** to show that the few-shot interactions indeed help to remove ambiguities and improve the performance. Furthermore, the

Where2Act-adaptation baseline helps substantiate the effectiveness of our proposed AAP module, while the **Ours-fps** baselines are designed to verify the usefulness of the proposed AIP module.

Besides, we compare to an ablated version of our method to verify the significance of iterative training between the AAP module and the AIP module.

- **Ours w/o iter**: an ablated version that trains the whole framework without the iterative training process.

Evaluation Metrics. Following Where2Act [18], we use the F-score, balancing the precision and recall, to evaluate the predictions of \mathcal{C}_{AAP} , and use the sample-successful rate (Sample-Succ) to evaluate the performance of \mathcal{C}_{AAP} and \mathcal{D}_{AAP} . To compute the sample successful rate, we apply the learned test-time strategy to fill \mathcal{I} and then use the extracted hidden information z as the conditional input to \mathcal{C}_{AAP} and \mathcal{D}_{AAP} . After that, we randomly select a point to interact from the group of points with the top-100 actionability scores $a_{p|z}^{AAP}$, sample 100 actions u_i 's at p , obtain $s_{u_i|z}^{AAP}$'s of these actions u_i 's predicted by \mathcal{C}_{AAP} , and then choose the action u_i with the highest $s_{u_i|z}^{AAP}$ to execute. We perform 10 interaction trials per test shape and report the final sample-succ rate as the percentage of sampling successful interactions in simulation.

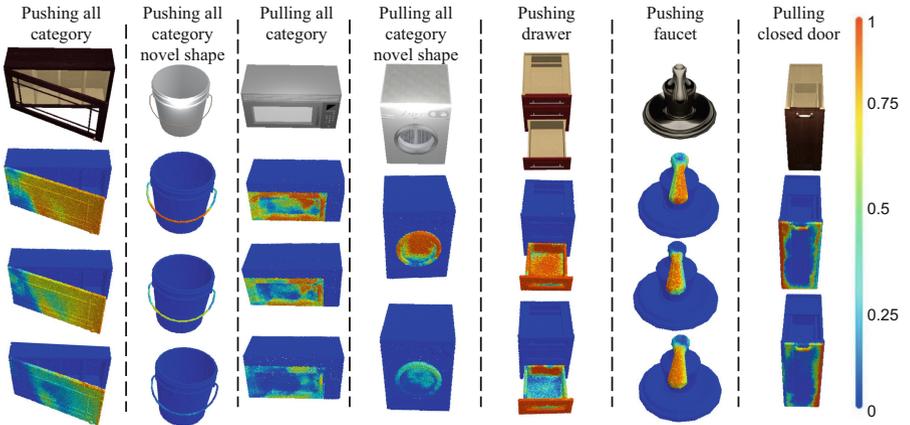


Fig. 4. Example results of adapted affordance predictions given by AAP module under different kinematic and dynamic parameters. The first five columns show the adapted affordance prediction conditioned on increasing joint friction (the first and second columns), part mass (the third column), and friction coefficient on object surface (the fourth and fifth columns). The last two columns respectively show the influence of different rotating directions (*i.e.*, joint limits) and joint axis locations.

5.3 Results and Analysis

Table 1 presents the quantitative comparisons against the baselines showing that our method achieves the best performance in most comparison entries. Specifically, compared to **Where2Act**, we observe that our method can improve the performance evidently with only 1 interaction. Also, the performance increases as the number of interactions increases in most cases. Compared to the **Where2Act-adaptation** baseline, our method with the proposed AAP module shows better performance, revealing that learning an adaptation network works better than using simple heuristics for adaptation. Finally, the superior performance against the **Ours-fps** baseline that use FPS sampled interaction trials further validate that our proposed AIP module is effective in strategically and iteratively picking interaction trials. Our method can generalize well to novel shapes and even shapes from unseen object categories through scores in test-cat.

Figure 4 shows example visualizations for our predicted affordance map posterior given interactions under different hidden kinematic or dynamic information (see the caption for more details explaining the different scenarios). In these figures, it is clear to see that our proposed method successfully adapts the affordance prediction conditioned on different hidden information. The affordance predictions within one shape share the same visual inputs but output different results, showing that our hidden embedding z contains certain information.

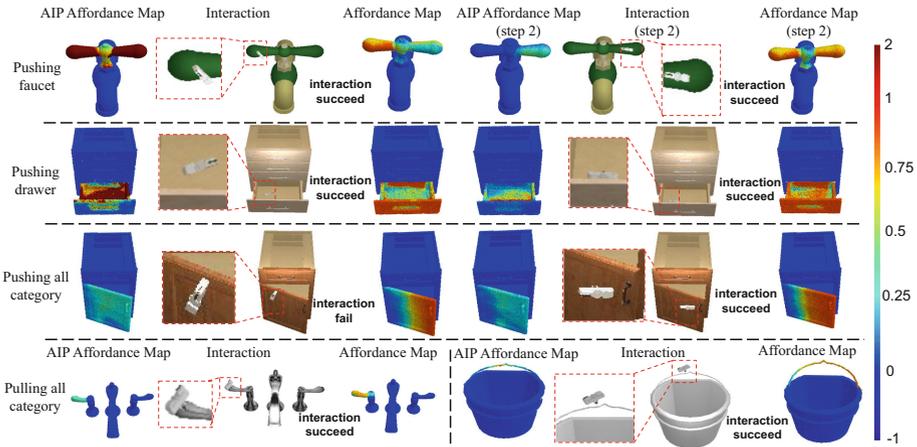


Fig. 5. Example results for the interactions proposed by the AIP module and the corresponding AIP affordance map predictions. In the first three rows, we show the initial and the second AIP affordance maps, the corresponding proposed interactions, and the posterior affordance map predictions. In the last row, we present two more examples that only one interaction is needed. From these results, our AIP module successfully proposes reasonable interactions for querying useful hidden information.

Figure 5 further shows some interaction proposals by our AIP module with its influence on the prediction of AAP affordance map and to the AIP affordance map itself. In the first row, for example, we see that the AIP affordance first proposes to interact at both sides of the faucet since it knows little about the hidden information but at the second timestep proposes the right side as it already learns the left side is actionable. Cases in the first and third rows demonstrate that the past few interactions will influence the selection of future interaction points, justifying the necessity of our recurrent structure for interaction selection. In the last row, we show cases only requiring one step to adapt.

Ablation Study. In Table 2, comparing against **Ours w/o iter** that trains the whole system without the interactive training process, we see that **Ours-final** achieves better results in most cases, which proves the effectiveness of the iterative training scheme. By iteratively alternating the training between the AAP module and the AIP module, the networks would be trained under the distribution of test-time interactions and thus achieve improved performance.

Real-World and Real-Robot Experiments. Finally, we perform real-world and real-robot experiments to show that our method can to some degree work beyond synthetic data. We use a Franka panda robot with a two-finger parallel gripper as the actuator to pull open a cabinet door. Figure 6 presents the results that our system proposes two interaction trials to inquire more information about this real-world cabinet and successfully learns to adapt to the posterior predictions. Please refer to the supplementary materials for a video better illustrating this example, more experiment settings, more example results, and more experiments with additional analysis.

Table 2. Ablation Study. We compare our method to an ablated version, where we remove the iteratively training process. It is clear to see that the iteratively training process helps our framework achieve better results in most cases.

		F-score (%)	Sample-Succ (%)
Pushing all (train cat.)	Ours w/o iter	71.21/72.64/73.16	30.67/31.62/32.56
	Ours-final	72.78/73.12/75.18	33.82/33.23/35.23
Pushing all (test cat.)	Ours w/o iter	77.24/77.33/77.17	31.03/33.89/ 38.83
	Ours-final	77.58/77.63/78.42	34.97/36.75/37.40
Pulling all (train cat.)	Ours w/o iter	41.19/42.10/42.81	6.67/7.22/8.33
	Ours-final	42.62/43.87/44.08	7.78/9.44/10.55
Pulling all (test cat.)	Ours w/o iter	48.31/48.28/50.50	5.65 /6.52/9.13
	Ours-final	49.51/50.00/51.33	5.21/ 7.39/10.45
Pulling closed door	Ours w/o iter	56.74/64.88/80.64	9.77/11.50/22.00
	Ours-final	57.83/65.60/79.65	10.86/11.57/22.14
Pushing faucet	Ours w/o iter	73.81/83.03/ 84.32	61.11/81.60/ 84.03
	Ours-final	77.42/83.06/83.83	65.90/81.66/82.14

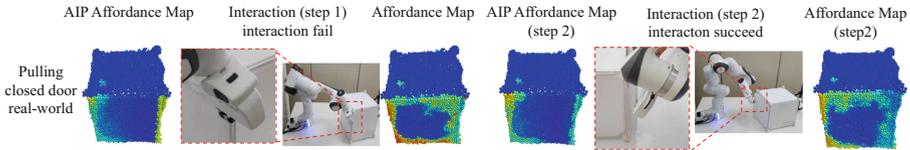


Fig. 6. Real-robot experiment on pulling open a closed door in the real world. We show the AIP affordance map predictions, the AIP proposed interactions, and the AAP posterior predictions, for two interaction trials. The results show that our work could reasonably generalize to real-world scenarios.

6 Conclusion

This work addresses a big limitation of previous works learning visual actionable affordance for manipulating 3D articulated objects – the hidden kinematic or dynamic uncertainties. We propose a novel framework AdaAfford that samples a few test-time interactions for fastly adapting to a more accurate affordance posterior prediction removing such ambiguities. Experimental results validate the effectiveness of our method compared to baseline approaches.

Limitations and Future Works. This work only considers two action types and 3D articulated objects. Future works may study more interaction and data types. Also, we only perform short-term interactions. Future works can investigate how to extend the framework for long-term manipulation trajectories. Future works shall work on considering the robot arm constraints.

Acknowledgements. National Natural Science Foundation of China -Youth Science Fund (No. 62006006). Leonidas and Kaichun were supported by the Toyota Research Institute (TRI) University 2.0 program, NSF grant IIS-1763268, a Vannevar Bush Faculty Fellowship, and a gift from the Amazon Research Awards program. The Toyota Research Institute University 2.0 program (Toyota Research Institute (“TRI”)) provided funds to assist the authors with their research but this article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity).

References

1. Agrawal, P., Nair, A., Abbeel, P., Malik, J., Levine, S.: Learning to poke by poking: experiential learning of intuitive physics. arXiv preprint [arXiv:1606.07419](https://arxiv.org/abs/1606.07419) (2016)
2. Chitta, S., Cohen, B., Likhachev, M.: Planning for autonomous door opening with a mobile manipulator. In: 2010 IEEE International Conference on Robotics and Automation, pp. 1799–1806. IEEE (2010)
3. Corona, E., Pumarola, A., Alenya, G., Moreno-Noguer, F., Rogez, G.: Ganhand: predicting human grasp affordances in multi-object scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5031–5041 (2020)
4. 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)

5. Farid, K., Sakr, N.: Few shot system identification for reinforcement learning. arXiv preprint [arXiv:2103.08850](https://arxiv.org/abs/2103.08850) (2021)
6. Ferreira, F., Shao, L., Asfour, T., Bohg, J.: Learning visual dynamics models of rigid objects using relational inductive biases. arXiv preprint [arXiv:1909.03749](https://arxiv.org/abs/1909.03749) (2019)
7. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: International Conference on Machine Learning, pp. 1126–1135. PMLR (2017)
8. Fouhey, D.F., Delaitre, V., Gupta, A., Efros, A.A., Laptev, I., Sivic, J.: People watching: human actions as a cue for single view geometry. In: Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., Schmid, C. (eds.) ECCV 2012. LNCS, vol. 7576, pp. 732–745. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-33715-4_53
9. Gibson, J.J.: The theory of affordances. *Hilldale USA* **1**(2), 67–82 (1977)
10. Janner, M., Levine, S., Freeman, W.T., Tenenbaum, J.B., Finn, C., Wu, J.: Reasoning about physical interactions with object-oriented prediction and planning. arXiv preprint [arXiv:1812.10972](https://arxiv.org/abs/1812.10972) (2018)
11. Jiang, Z., Zhu, Y., Svetlik, M., Fang, K., Zhu, Y.: Synergies between affordance and geometry: 6-DoF grasp detection via implicit representations. In: Proceedings of Robotics: Science and Systems (RSS) (2021)
12. Kjellström, H., Romero, J., Kragić, D.: Visual object-action recognition: inferring object affordances from human demonstration. *Comput. Vis. Image Underst.* **115**(1), 81–90 (2011)
13. Kovic, M., Kragić, D., Bohg, J.: Learning task-oriented grasping from human activity datasets. *IEEE Robot. Autom. Lett.* **5**(2), 3352–3359 (2020)
14. Kumar, K.N., Essa, I., Ha, S., Liu, C.K.: Estimating mass distribution of articulated objects using non-prehensile manipulation. arXiv preprint [arXiv:1907.03964](https://arxiv.org/abs/1907.03964) (2019)
15. Lenz, I., Lee, H., Saxena, A.: Deep learning for detecting robotic grasps. *Int. J. Robot. Res.* **34**(4–5), 705–724 (2015)
16. Li, X., Liu, S., Kim, K., Wang, X., Yang, M.H., Kautz, J.: Putting humans in a scene: learning affordance in 3D indoor environments. In: IEEE Conference on Computer Vision and Pattern Recognition (2019)
17. Mandikal, P., Grauman, K.: Learning dexterous grasping with object-centric visual affordances. In: IEEE International Conference on Robotics and Automation (ICRA) (2021)
18. Mo, K., Guibas, L.J., Mukadam, M., Gupta, A., Tulsiani, S.: Where2act: from pixels to actions for articulated 3D objects. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 6813–6823, October 2021
19. Mo, K., Qin, Y., Xiang, F., Su, H., Guibas, L.: O2O-afford: annotation-free large-scale object-object affordance learning. In: Conference on Robot Learning (CoRL) (2021)
20. Mo, K., et al.: PartNet: a large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2019
21. Nagarajan, T., Grauman, K.: Learning affordance landscapes for interaction exploration in 3D environments. In: NeurIPS (2020)
22. NVIDIA. [Nvidia.physx](https://github.com/NVIDIA/physx)

23. Peterson, L., Austin, D., Kragic, D.: High-level control of a mobile manipulator for door opening. In: Proceedings. 2000 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2000) (Cat. No. 00CH37113), vol. 3, pp. 2333–2338. IEEE (2000)
24. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: deep hierarchical feature learning on point sets in a metric space. arXiv preprint [arXiv:1706.02413](https://arxiv.org/abs/1706.02413) (2017)
25. Qi, W., Mullapudi, R.T., Gupta, S., Ramanan, D.: Learning to move with affordance maps. arXiv preprint [arXiv:2001.02364](https://arxiv.org/abs/2001.02364) (2020)
26. Qin, Y., Chen, R., Zhu, H., Song, M., Xu, J., Su, H.: S4G: amodal single-view single-shot se (3) grasp detection in cluttered scenes. In: Conference on Robot Learning, pp. 53–65. PMLR (2020)
27. Qin, Z., Fang, K., Zhu, Y., Fei-Fei, L., Savarese, S.: Keto: learning keypoint representations for tool manipulation. In: 2020 IEEE International Conference on Robotics and Automation (ICRA), pp. 7278–7285. IEEE (2020)
28. Rakelly, K., Zhou, A., Finn, C., Levine, S., Quillen, D.: Efficient off-policy meta-reinforcement learning via probabilistic context variables. In: International Conference on Machine Learning, pp. 5331–5340. PMLR (2019)
29. Redmon, J., Angelova, A.: Real-time grasp detection using convolutional neural networks. In: 2015 IEEE International Conference on Robotics and Automation (ICRA), pp. 1316–1322. IEEE (2015)
30. Schmidt, T., Newcombe, R.A., Fox, D.: Dart: dense articulated real-time tracking. In: Robotics: Science and Systems, Berkeley, CA, vol. 2 (2014)
31. Sun, Yu., Ren, S., Lin, Y.: Object-object interaction affordance learning. *Robot. Autom. Syst.* **62**(4), 487–496 (2014)
32. Tzionas, D., Gall, J.: Reconstructing articulated rigged models from RGB-D videos. In: Hua, G., Jégou, H. (eds.) ECCV 2016. LNCS, vol. 9915, pp. 620–633. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-49409-8_53
33. Urakami, Y., Hodgkinson, A., Carlin, C., Leu, R., Rigazio, L., Abbeel, P.: Doorgym: a scalable door opening environment and baseline agent. In: Deep RL workshop at NeurIPS 2019 (2019)
34. Wang, X., Zhou, B., Shi, Y., Chen, X., Zhao, Q., Xu, K.: Shape2motion: joint analysis of motion parts and attributes from 3D shapes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8876–8884 (2019)
35. Weng, Y., et al.: Captra: category-level pose tracking for rigid and articulated objects from point clouds. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV), pp. 13209–13218, October 2021
36. Wu, R., et al.: VAT-mart: learning visual action trajectory proposals for manipulating 3D ARTiculated objects. In: International Conference on Learning Representations (2022)
37. Xiang, F., et al.: SAPIEN: a simulated part-based interactive environment. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2020
38. Xu, Z., Wu, J., Zeng, A., Tenenbaum, J.B., Song, S.: Densephysnet: learning dense physical object representations via multi-step dynamic interactions. arXiv preprint [arXiv:1906.03853](https://arxiv.org/abs/1906.03853) (2019)
39. Xu, Z., He, Z., Song, S.: UMPNet: universal manipulation policy network for articulated objects. *IEEE Robot. Autom. Lett.* (2022)
40. Yan, Z., et al.: RPM-NET: recurrent prediction of motion and parts from point cloud. *ACM Trans. Graph.* **38**(6), Article 240 (2019)

41. Yang, L., Zhan, X., Li, K., Xu, W., Li, J., Lu, C.: CPF: learning a contact potential field to model the hand-object interaction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 11097–11106 (2021)
42. Yi, L., Huang, H., Liu, D., Kalogerakis, E., Su, H., Guibas, L.: Deep part induction from articulated object pairs. *ACM Trans. Graph.* **37**(6) (2018)
43. Yu, W., Tan, J., Liu, C.K., Turk, G.: Preparing for the unknown: learning a universal policy with online system identification. arXiv preprint [arXiv:1702.02453](https://arxiv.org/abs/1702.02453) (2017)
44. Zhao, T.Z., Nagabandi, A., Rakelly, K., Finn, C., Levine, S.: Meld: meta-reinforcement learning from images via latent state models. arXiv preprint [arXiv:2010.13957](https://arxiv.org/abs/2010.13957) (2020)
45. Zhou, W., Pinto, L., Gupta, A.: Environment probing interaction policies. In: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, 6–9 May 2019. OpenReview.net (2019)
46. Zhu, Y., Zhao, Y., Chun Zhu, S.: Understanding tools: task-oriented object modeling, learning and recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2855–2864 (2015)