

# Learning to Plan with Pointcloud Affordances for General-Purpose Dexterous Manipulation

Anthony Simeonov<sup>1</sup>, Yilun Du<sup>1</sup>, Beomjoon Kim<sup>1</sup>, Francois R. Hogan<sup>2</sup>, Alberto Rodriguez<sup>1</sup>, and Pulkit Agrawal<sup>1</sup>

**Abstract**—We aim to enable a robot to solve arbitrary manipulation tasks using dexterous manipulation primitives. This suggests the use of techniques from task-and-motion planning, which can sequence primitives by performing multistep reasoning and forward simulation. At the same time, we aim to solve such tasks when presented with objects of different shapes and sizes, suggesting the use of learned perceptual representation which enable generalization across geometries. However, it has remained challenging to incorporate these representations into systems that perform accurate forward simulation, a necessary component for reasoning toward long-horizon goals. We propose a framework, that utilizes deep generative models and segmented object pointclouds, that enables multistep planning using dexterous primitive manipulation skills in tasks involving a variety of object shapes and sizes. Our learned models, taking the form of biased sampling distributions, provide gains in planning efficiency over a manually designed baseline when integrated with a sampling-based planner. We also contribute a set of novel design choices in this framework which provide benefits in generalization and sample quality.

## I. INTRODUCTION

Consider the task depicted in Figure 1. A two-arm robot must use its palms to move the red box from its initial configuration to the green goal pose, which is on the opposite side of the table and on a different face of the object. This can be imagined as a proxy for the real-world task of moving, for example, a book from an arbitrary initial tabletop pose into a specified upright pose in the corner of a shelf.

We assume that the robot has access to a variety of parameterized primitive skill behaviors that can be combined to solve the task, such as pulling, to translate and change the yaw angle of the object, and grasping, to flip it onto a different face. Parameter values for each skill must be chosen that determine *how* the skill is executed. Finding a task solution amounts to searching for a suitable sequence of skill types and corresponding parameter values (i.e. *where* to grasp and *how* to flip) for a particular start and goal pose of the object.

One way to tackle this problem is by using a task-and-motion planning (TAMP) algorithm [7, 2, 5, 12], which usually incorporates construction of a search tree, sampling from distributions over skill parameters, and iterative simulation of many skill sequences until one that reaches the goal is found. One of the major drawbacks of most of these TAMP algorithms is that they assume access to the pose and shape of the objects in the scene, and use manually-designed samplers that use this information for skill parameters.

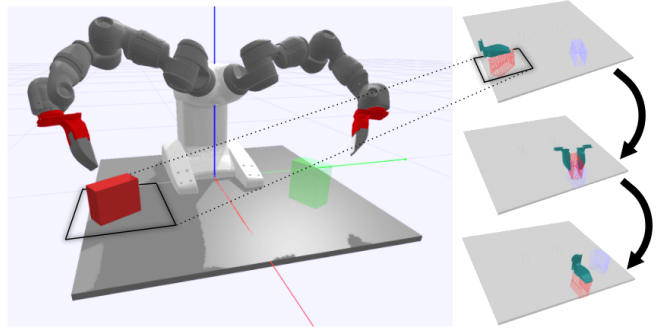


Figure 1. Our framework uses conditional generative models and dexterous primitives to imagine and execute multi-step manipulation plans for solving arbitrary start/goal manipulation tasks, when presented with only a pointcloud observation of the object.

For robots to operate in unforeseen environments however, they must be able to manipulate objects whose shape and pose are not known apriori. For example, consider again the problem in Fig. 1. The robot may need to put away many different sized books, but not know in advance their specific size and shape. This requires a sampler that can directly operate on sensory input, which is difficult to hand-design.

Based on this observation, we propose a framework that learns skill parameter samplers from prior experience using the primitive skills [8], that directly operate on sensory observations. While our framework can be used within any TAMP algorithm, in this work, we assume that a *plan skeleton* [9], which is a sequence of primitive types, such as pull-grasp-push, is given. Our objective is to find the continuous parameters of the plan skeleton using the learned samplers.

One key challenge in designing such a framework is learning the representation of an object represented by unstructured sensor data that enables both long-horizon reasoning and generalization across object shapes. There has been much work on learning representations of RGB-D images that are tailored to distinct action types, such as top-down grasping [15, 14, 16], suction [15] and pushing [1, 14]. These learning-based techniques for encoding sensory observations demonstrate strong generalization to diverse real world objects and are thus applicable for our similar goals of object generalization.

However, it has remained difficult to integrate these representations into systems that can solve problems involving longer planning horizons, in part because directly learning highly accurate forward models in this space of sensory data

<sup>1</sup>Massachusetts Institute of Technology. <sup>2</sup>Samsung Electronics. Correspondence to: Anthony Simeonov <asimeono@mit.edu>

is extremely difficult. We instead propose to use segmented pointclouds as an object representation that provides more structure than pixel-based representations but is still flexible enough to afford generalization across object shapes.

To support long-horizon planning, we propose a system design where we learn two different samplers: a subgoal sampler for predicting a reachable rigid body transform that can be used to forward simulate and imagine future pointcloud observations, and a contact sampler that generates end-effector poses suitable for achieving the predicted subgoal. The samplers are designed to learn and exploit the correlation between contacts and subgoals. To support generalization across objects, we use recent advancements in neural network architectures that can operate on pointcloud data, such as PointNet++ [10] and Graph Attention Networks [13], in a conditional generative modeling scheme where the samplers are represented as neural networks.

We validate our approach in a simulated domain where a dual arm system, equipped with end effector palms, is tasked with solving a large distribution of object manipulation tasks. Our method provides significant advantages in planning efficiency and prediction quality over a manually designed baseline that utilizes privileged knowledge about the shape of the objects and the task.

## II. PROBLEM DEFINITION

### A. Problem Setup

We define a robot manipulation primitive *skill*  $\pi$  as a function that takes as an input a set of parameters,  $\Theta$ , and outputs a robot joint trajectory. For instance, consider the right-hand pulling skill, denoted  $\pi_{\text{PULLRH}}$ . Its parameters, denoted  $\Theta_{\text{PULLRH}}$ , are the pose of the right palm,  $T_R^p \in SE(3)$ , expressed in the world coordinate frame, and the desired rigid body transformation of the object,  $T^o \in SE(3)$ . Given  $T_R^p$  and  $T^o$ , the primitive skill outputs a sequence of right arm configurations  $(q_0, \dots, q_M)$ ,  $q \in \mathbb{R}^d$ , for pulling the object. The skill is feasible if  $(q_0, \dots, q_M)$  is collision-free, does not have large joint velocities, and does not go through singularities. If  $T_R^p$  leads to a compatible contact configuration between the end effector and the object at its initial pose, then when  $(q_0, \dots, q_M)$  is followed,  $T^o$  will be applied to the object.

We consider general primitives of this type, which take as parameters an initial world frame pose of one or both robot palms  $T_{R,L}^p$  and a desired object transformation  $T^o$ ,

$$\pi : T_R^p \times T_L^p \times T^o \longrightarrow \mathcal{Q}_{R,L}^*$$

where  $\mathcal{Q}_{R,L}^*$  denotes the space of right and left arm configuration sequences. In particular, we utilize the primitive planning scheme developed in [6], although other primitives parameterized by an initial contact configuration and a desired motion to apply are equally applicable to the presented framework.

### B. Planning Problem Definition

Assume we are given a desired rigid body transformation to be applied to the object,  $T_{des}^o \in SE(3)$  along with a *plan*

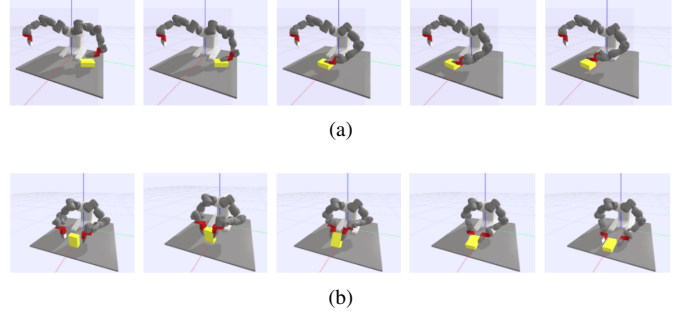


Figure 2. Snapshots during the execution of a (a) pull skill and (b) grasp skill on the simulated YuMi robot in PyBullet.

*skeleton*  $PS$  that defines the high level sequence of primitives that should be used to complete the task:

$$PS = \pi_0(\Theta_0) \longrightarrow \pi_1(\Theta_1) \longrightarrow \dots \longrightarrow \pi_K(\Theta_K)$$

For plan skeleton  $PS$ , we denote its space of parameters as  $\Theta_{PS}$ . The robot observes the world using a set of RGB-D sensors which provide a segmented point cloud observation of the environment  $X \in \mathbb{R}^{N \times 3}$ . Given a set of primitive skills, a plan skeleton  $PS$ , a desired object transformation  $T_{des}^o$ , and a pointcloud observation  $X$  of the object at some initial configuration, our objective is to find the parameters of the given plan skeleton,  $\theta_{PS} \in \Theta_{PS}$ , that achieves  $T_{des}^o$ .

### C. Learning Problem Definition

Given a distribution from which skill parameters can be sampled, an RRT-style algorithm (described in detail in Section III-D) can be used to solve the planning problem described in Section II-B. This involves sampling intermediate values of  $T^o$ , referred to as *subgoals*, and corresponding end effector poses  $T^p$  that can reach these subgoals. The intuition behind this approach is the following: A generative model over subgoals can act a forward model, since sampled values of  $T^o$  can be used to transform an observed  $X_t$  into an imagined  $\hat{X}_{t+1}$ . This process can repeat iteratively, enabling multiple steps of forward simulation. At the same time, a generative model over contact poses acts as a biased action sampler, that produces actions that are likely to be feasible for the particular object being represented by  $X$ . When combined, pairs of subgoals and contact poses can be provided as input to  $\pi(\Theta)$  to determine the feasibility of the resulting motion.

We are given an experience dataset  $D$  consisting of successful single-step skill executions. These can be considered as a set of demonstrations of solving planning problems with  $PS$  of length *one*, where the subgoal parameter  $T^o$  equals  $T_{des}^o$ .

$$D = \{(X^{(i)}, PS^{(i)}, T_{des}^{o(i)}, \theta_{PS}^{(i)})\}_{i=1}^n$$

Using this data, we aim to train *conditional generative skill models*,  $p_{\phi, \text{SKILL}}(\cdot|X)$ , taking the form of a deep neural network parameterized by  $\phi$ , that at test time produces a distribution over high likelihood skill parameters when provided with a new pointcloud observation.

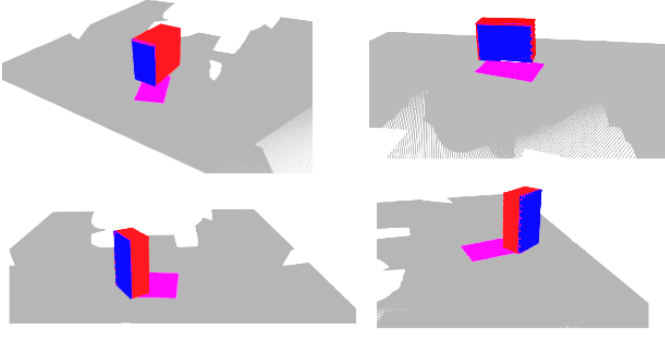


Figure 3. Labeled training data for the task of predicting object-table contact segmentation masks. The blue points are those on the object pointcloud that end up in contact with the table after  $\pi_{\text{GRASPREORIENT}}$  skill execution (pink points on the table are shown for visualization of where the blue points end up). Given the masked blue points and the table pointcloud, registration can be used to solve for stable subgoal  $T^o$ .

### III. METHODOLOGY

#### A. Learning Approach

We leverage conditional variational auto-encoders (CVAE) [11] to learn  $p_{\phi, \text{SKILL}}$  using dataset  $D$ . This approach raises multiple open questions in the design of the system architecture and skill parameter representation, and these questions are discussed in the following sections.

#### B. Learning the joint subgoal-contact distribution

One naive application of this framework to our problem of sampling  $T^p$  and  $T^o$  is to train separate CVAEs, effectively modeling  $T^o$  and  $T^p$  as independent random variables. Another option is to model their dependence via a conditional distribution, i.e. learn one CVAE that predicts subgoals conditioned on pointclouds, and learn a second CVAE that predicts contact configurations conditioned on pointclouds and subgoals. Instead, our approach is to directly model the *joint* distribution via a single latent variable, i.e. *simultaneously* predict both subgoals and contact configurations conditioned on pointclouds.

#### C. Generalizable $SE(3)$ Subgoal Representation

Subgoals for skills such as  $\pi_{\text{GRASPREORIENT}}$  (Fig. 2b) must be represented in full  $SE(3)$ . For generalization to novel geometries and object states, it is unclear whether directly predicting an  $SE(3)$  transformation is the most useful approach, as we require the ability to predict  $T^o$  that will lead to a reachable and passively stable object pose that doesn't change when the robot breaks contact and moves to the next step in  $PS$ . Directly predicting  $T^o$  is not well suited to meet this requirement, and can lead to transforming the pointcloud into unstable or unreachable configurations (i.e. floating in air, penetrating the table, tilted on an edge, etc.). We instead propose to have the model predict a binary segmentation mask directly with respect to the observed pointcloud, representing the set of points that *should end in contact with the table*, and then use these points in a registration routine to solve for  $T^o$ .

Fig. 3 shows a visual depiction of this representation on a few cuboidal objects.

#### D. Generalizable Multistep Planning With Pointclouds

We integrate our learned generative models with a sampling-based planning framework inspired by RRT. For the first step in  $PS$ , the initial pointcloud  $X_0$  is used to sample parameters  $\theta_0 \sim p_{\phi, \text{SKILL}}(\cdot | X_0)$  for that corresponding skill. If the instantiated skill  $\pi_0(\theta_0)$  is feasible,  $\theta_0$  are saved for that position in  $PS$ . Additionally, the initial pointcloud  $X_0$  is transformed via the sampled  $T_0^o$  into a new configuration  $X_1$  that can be used as the conditioning variable when sampling the *next* skill in the skeleton, and this process is repeated. For the final step in  $PS$ , the required unknown  $T_K^o$  can be solved for based on the required transformation-to-go to solve the task as  $T_K^o = T_{des}^o \prod_{i=0}^{K-1} T_i^o$ , since we know the determined returned subgoals  $T_0^o, \dots, T_K^o$  should follow  $T_{des}^o = \prod_{i=0}^K T_{K-i}^o$ .

### IV. EXPERIMENTS AND RESULTS

Our experiments are designed to validate our primary claim that our learned sampler performs better than a hand-designed sampler when integrated with our multistep planning framework. We also conduct experiments to answer whether our choice of joint subgoal-contact distribution modeling and segmentation mask subgoal representation lead to better generated sample quality and generalization to unseen pointclouds. In the current work we focus on the primitive  $\pi_{\text{GRASPREORIENT}}$ .

#### A. Experimental Setup

The framework is implemented on an ABB YuMi robot simulated in PyBullet [3]. GelSlim [4] end effectors are used, similarly to as in [6], without the tactile sensing simulated. The robot operates in a table-top environment with RGB-D cameras located at the four corners all focused at the same point in the front of the robot. We use the built-in segmentation mask capability in PyBullet to obtain segmented pointclouds of the object from the simulated depth images.

#### B. Training Details and Network Architecture

We compare PointNet++ [10] and Graph Attention Networks (GAT) [13] for the encoder and decoder in our CVAE, as these architectures have shown success in learning useful embeddings for downstream tasks involving graph-structured data such as pointclouds. The encoder is trained to map the data in  $D$  to a latent conditional distribution, constrained by a KL-divergence loss to resemble a unit Gaussian. Sampled latents are concatenated as an additional feature to the points in  $X$ , which the decoder uses to reconstruct the data.

We generate a distribution of cuboids of varying dimensions. The model is trained on a small subset of them and tested on cuboids with dimensions never seen during training. The data generation procedure relies on the simulation of the primitive skills in a scenario where 3D object models are available, allowing computation of stable poses, collision checking between the robot and the environment, and using a rejection sampling scheme to randomly sample skill parameters to find

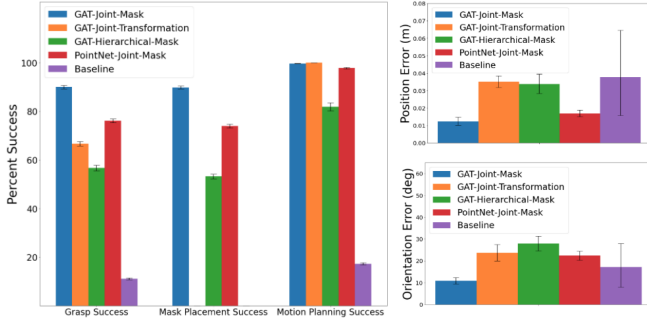


Figure 4. Single Step Grasp+Reorient Results: Modeling the joint vs. the conditional distribution leads to substantially better generalization, and our segmentation mask subgoal representation leads to significantly lower position and orientation errors.

ones that are feasible. When feasible skills are found, they are executed and the parameters are added to the dataset.

### C. Evaluation: Single-step Grasp + Reorient

To demonstrate the value of the novel aspects of our approach, we conduct a set of single-step trials using variants of the learned sampler. For 20 unseen cuboid objects of different dimensions, 60 different start poses in a region near the front of the robot are sampled per object. From the pointcloud observation of the object in each start state, the learned model is tasked with producing its own  $T^o$  and  $T^p$ .

We compare both subgoal representations, one which involves directly predicting  $T^o$ , along with our novel approach that involves predicting a binary per-point segmentation mask, and using the masked points in a registration step to solve for  $T^o$ . We use ICP to solve for the registration, and initialize the registration with a pure forward  $\frac{\pi}{2}$  pitch about the object body frame. We also implement a manually designed baseline using plane segmentation and antipodal point heuristics based on the privileged knowledge that we are operating with cuboids. Note that our method never assumes anything about the global geometry of the manipulated objects.

For each start state, we measure the following

- Grasp stability: Did the object miss making contact, or drop the object during execution?
- Mask placement success: Did the set of predicted table-object contact points end up on the table?
- Motion planning success: Given a budget of 15 samples, could the motion planner find a feasible plan?
- Pose error: For feasible plans that are executed, how close was the executed object transformation to the sampled  $T^o$

Fig. 4 shows the GAT-joint-mask model performs best.

### D. Evaluation: Multi-step Manipulation

To validate the benefits provided by our learned sampler, we conduct multi-step planning experiments using unseen objects and measure fixed-budget planning success rate and returned plan quality for a variety of plan skeletons. The framework in III-D is implemented using the best performing

Table I  
MULTISTEP PLANNING RESULTS

Planning Success Rate	PG	GP	PGP
Learned	0.78	0.83	0.79
Uniform	0.19	0.19	0.05
Sticking Contact Success Rate	PG	GP	PGP
Learned	0.79	0.93	0.74
Uniform	0.75	0.71	0.50

Multistep planning success rate with a fixed 5-minute timeout using our learned skill samplers and a manually designed uniform sampler. Planning success rate indicates percent of planning problems where the planner found a solution before timing out. Sticking contact success rate indicates percent of plans that were executed where contact was maintained with the object for the duration of each step where sticking contact is assumed. PG: Pull  $\rightarrow$  Grasp, GP: Grasp  $\rightarrow$  Pull, PGP: Pull  $\rightarrow$  Grasp  $\rightarrow$  Pull

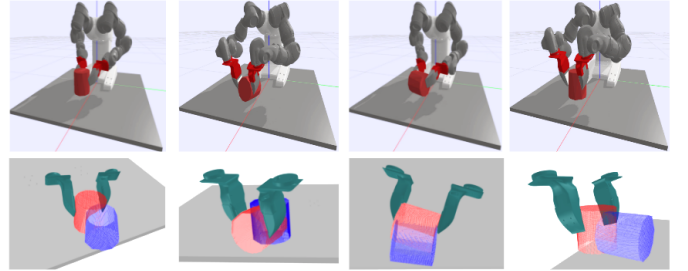


Figure 5. Geometry Generalization: The best performing GAT-Joint-Mask  $p_{\phi, \text{GRASP REORIENT}}$  model was trained only on cuboids, and tested on cylinders in a small set of configurations. Top row: simulator execution of predictions. Bottom row: visualization of subgoal (blue pointclouds) and contact (green palm) predictions. The predictions were qualitatively observed to be feasible, indicating the model has somewhat learned to generalize to global geometry outside the training distribution.

$p_{\phi, \text{GRASP REORIENT}}$  model. The same heuristic baseline is used for comparison, and a similar heuristic sampler for  $\pi_{\text{PULL RH}}$  (Fig. 2a) is used in both cases.

Planning success rate quantifies the fraction of trials where a plan is found before timeout. We also track whether the robot maintains contact during open loop execution to quantify the quality of the plans that are returned. Table I results indicate that the learned sampler provides a large benefit in planning efficiency and returned plan quality.

### E. Qualitative Results: Geometry Generalization

Fig. 5 shows predictions made by models trained only on cuboids, and tested on cylinders in different configurations. As shown, the predictions for both  $T^o$  and  $T^p$  are quite sensible and were able to be executed successfully in the simulator.

## V. CONCLUSION

We presented a method based on deep generative modeling applied to segmented pointclouds to enable general-purpose multistep sampling-based planning using dexterous primitive manipulation skills. Our method enables planning efficiency gains over a manually designed baseline sampler, while simultaneously allowing generalization to unseen objects by utilizing a perception-driven object representation.



## REFERENCES

- [1] Pulkit Agrawal, Ashvin V Nair, Pieter Abbeel, Jitendra Malik, and Sergey Levine. Learning to poke by poking: Experiential learning of intuitive physics. In *Advances in neural information processing systems*, pages 5074–5082, 2016.
- [2] S. Cambon, R. Alami, and F. Gravot. A hybrid approach to intricate motion, manipulation, and task planning. *International Journal of Robotics Research*, 2009.
- [3] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. *GitHub repository*, 2016.
- [4] Elliott Donlon, Siyuan Dong, Melody Liu, Jianhua Li, Edward Adelson, and Alberto Rodriguez. Gelslim: A high-resolution, compact, robust, and calibrated tactile-sensing finger. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1927–1934. IEEE, 2018.
- [5] C. R. Garrett, T. Lozano-Pérez, and L. P. Kaelbling. Ffrob: Leveraging symbolic planning for efficient task and motion planning. *International Journal of Robotics Research*, 2014.
- [6] Francois R. Hogan, Jose Ballester, Siyuan Dong, and Alberto Rodriguez. Tactile dexterity: Manipulation primitives with tactile feedback, 2020.
- [7] Leslie Pack Kaelbling and Tomas Lozano-Perez. Hierarchical task and motion planning in the now. In *IEEE Conference on Robotics and Automation (ICRA)*, 2011. URL <http://people.csail.mit.edu/lpk/papers/hpnICRA11Final.pdf>. Finalist, Best Manipulation Paper Award.
- [8] Beomjoon Kim, Leslie Pack Kaelbling, and Tomas Lozano-Perez. Guiding search in continuous state-action spaces by learning an action sampler from off-target search experience. In *Proceedings of the 32th AAAI Conference on Artificial Intelligence (AAAI)*. To appear. AAAI Press, 2018. URL <http://lis.csail.mit.edu/pubs/kim-aaai18.pdf>.
- [9] Tomás Lozano-Pérez and Leslie Pack Kaelbling. A constraint-based method for solving sequential manipulation planning problems. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3684–3691. IEEE, 2014.
- [10] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in neural information processing systems*, pages 5099–5108, 2017.
- [11] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In *Advances in neural information processing systems*, pages 3483–3491, 2015.
- [12] Marc Toussaint. Logic-geometric programming: An optimization-based approach to combined task and motion planning. *International Joint Conference on Artificial Intelligence*, 2015.
- [13] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- [14] Andy Zeng, Shuran Song, Stefan Welker, Johnny Lee, Alberto Rodriguez, and Thomas Funkhouser. Learning synergies between pushing and grasping with self-supervised deep reinforcement learning. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4238–4245. IEEE, 2018.
- [15] Andy Zeng, Shuran Song, Kuan-Ting Yu, Elliott Donlon, Francois R Hogan, Maria Bauza, Daolin Ma, Orion Taylor, Melody Liu, Eudald 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)*, pages 1–8. IEEE, 2018.
- [16] Andy Zeng, Shuran Song, Johnny Lee, Alberto Rodriguez, and Thomas Funkhouser. Tossingbot: Learning to throw arbitrary objects with residual physics. *arXiv preprint arXiv:1903.11239*, 2019.