

Predicting Motion Plans for Articulating Everyday Objects

Arjun Gupta Max E. Shepherd Saurabh Gupta

Abstract—Mobile manipulation tasks such as opening a door, pulling open a drawer, or lifting a toilet lid require constrained motion of the end-effector under environmental and task constraints. This, coupled with partial information in novel environments, makes it challenging to employ classical motion planning approaches at test time. Our key insight is to cast it as a learning problem to leverage past experience of solving similar planning problems to directly predict motion plans for mobile manipulation tasks in novel situations at test time. To enable this, we develop a simulator, ArtObjSim, that simulates articulated objects placed in real scenes. We then introduce SeqIK+ θ_0 , a fast and flexible representation for motion plans. Finally, we learn models that use SeqIK+ θ_0 to quickly predict motion plans for articulating novel objects at test time. Experimental evaluation shows improved speed and accuracy at generating motion plans than pure search-based methods and pure learning methods.

I. INTRODUCTION

As humans, when faced with everyday articulated objects as shown in Figure 1, we draw upon our vast past experience to successfully articulate them. We know to stand on the side as we pull open a oven, and where to lean on a door to push it open. Very rarely do we pull open a door onto our feet, or bump into the toilet while lifting a toilet seat. In this paper, we develop techniques that enable robots to similarly use past experience to *mine* and quickly *predict* strategies for articulating everyday objects in cluttered real environments.

Current work on articulating objects casts it as a motion planning problem: given a full scan of the environment, find a robot joint trajectory that leads the end-effector to track the trajectory that the grasp-point on the object should follow. This suffers from both a high-sensing cost and a high-planning cost. Building a full articulable 3D reconstruction of the environment for collision checking and planning is expensive and time consuming. At the same time, finding paths that conform to tight constraints on the end-effector trajectory while not colliding with self or surrounding obstacles or the articulating object is computationally hard. States that adhere to the given constraint form a measure zero set among the set of all states. This creates issues for sampling-based motion planners which can fail to sample states that satisfy the constraint, or must incur computation cost to project states to the constraint manifold [4], [20].

Rather than re-solving, from scratch, how to open a door every time we encounter one, our proposal is to build a repertoire of strategies based on past experience. This replaces the search in the high-dimensional motion plan space with the



Fig. 1: Household robots need to articulate everyday objects (e.g. pull open drawers, swing open cupboards, lift toilet seats). Such articulation involves applying forces onto the environment while maintaining relevant contact, such as with the drawer handle as we pull it open. This requires reasoning about the feasibility of the entire trajectory (i.e. points along the trajectory should not just be reachable, but it must be possible to continuously go from one point to the next). This paper develops datasets and techniques for learning models that can predict motion plans for such constrained motion planning problems with low sensing and planning costs.

much simpler problem of selecting from a small family of good strategies, leading to gains in efficiency. Furthermore, this simpler search can be driven by whatever observation is readily available from on-board sensors through the use of machine learning. Our experiments demonstrate the effectiveness of casting this as a learning problem. Given a single RGB-D observation of an articulated object in cluttered real world scenes and associated end-effector pose trajectory to track, we can output motion plans that track the end-effector trajectory to within 0.01m and 0.01rad error with just a few inverse kinematic calls. This, by far, outperforms the constrained motion planning implementation for the projected state space method from the OMPL library [20], [44] which fails to find any motion plans with less than 0.01m and 0.01rad tracking error even when given 15 minutes of planning time. Our impressive performance is enabled by the following three innovations.

First, we construct, ArtObjSim, a lightweight kinematic

Authors are with the University of Illinois, Urbana-Champaign. EMail: {arjung2, maxes2, saurabhg}@illinois.edu. Project website with more details: <https://arjung128.github.io/mpao>.

simulator for everyday articulated objects placed in real scenes. Crucially, this simulator is derived from scans of *real-world* environments (from HM3D dataset [37]). This retains the appearance and the cluttered environmental context of the articulated objects. The simulator not only provides the experience to build the repertoire of strategies, but also serves as the first of its kind benchmark for evaluating motion plans for articulating objects in real environments. ArtObjSim consists of 3758 articulated object instances across 4 articulation types (prismatic *e.g.* drawers, vertical hinge *e.g.* cabinets, horizontal up-hinge *e.g.* toilet lids, horizontal down-hinge *e.g.* dishwashers) across 10 object categories in 97 scenes.

Second, rather than predicting a motion plans, that must conform to tight task constraints and are hence hard to directly predict, we instead predict a *strategy* that can be efficiently decoded into a motion plan using the articulation geometry. Our decoding process consists of *synchronously* solving inverse kinematics (IK) problems for end-effector waypoints sampled along the given end-effector trajectory. This synchronization is done by warm starting IK for the t^{th} time-step using solution from the $(t - 1)^{\text{th}}$ time-step. We call this decoding process Sequential Inverse Kinematics or SeqIK. By directly optimizing to reduce end-effector pose error, SeqIK leads to low tracking errors. The initialization for the first time step, θ_0 , serves as the strategy. Changing θ_0 changes the strategy and generates a different motion plan. We find that this representation, SeqIK+ θ_0 , is fast (motion plans can be quickly decoded) and flexible (with the right θ_0 it can produce high-quality motion plans for diverse objects in diverse situations).

Not all initializations would work well for all situations. Some might not be able to track the end-effector accurately enough, some may lead to collisions, and others yet might violate the task constraint when joint angles are interpolated for smooth execution. Thus, we need to find good initializations for SeqIK+ θ_0 at test time. Our third innovation, the use of a convolutional neural network to predict good initializations for SeqIK+ θ_0 (or equivalently, good strategies) from RGB image observations, speeds up test time inference. We train this model on a dataset of object images labeled with good initializations, as generated using our proposed ArtObjSim simulator. We are able to find good solutions with only a few IK calls. This is much faster than sampling-based planning at test time which would make tens of thousands of IK calls to project sampled states to the constraint set. We also show that our method can work with predicted end-effector waypoints. Collected dataset, ArtObjSim, and code are available on the project website: <https://arjung128.github.io/mpao/>.

II. RELATED WORK

Motion planning under constraints [4], [20] has been used to tackle object articulation problems, *e.g.* [5], [6], [8], [29], [36], [38], [46] among numerous other works. Researchers have tackled many aspects: design of task-space regions for expressing constraints on end-effectors [5], planning for base and arm motion separately [29], considering whole-body manipulation [6], reasoning about good locations to position

the base through inverse reachability maps [46], and even casting it instead as a trajectory optimization or optimal control problem [10], [30], [35], [42]. All these approaches solve a new object articulation problem, from scratch, every time they encounter one. Consequently, they incur a high sensing and planning costs. Different from these works, our interest is in techniques to leverage experience with similar articulation problems in the past to quickly predict motion plans with low sensing and planning cost. Online system identification approaches [17], [19], [33], [34] that adapt plans using feedback have also been studied.

Perception of articulated objects. A body of work [1]–[3], [18], [22], [27], [31], [32], [39], [47]–[49], [52] has tackled the perception of *articulation* geometry for articulated objects. Given raw sensory input (RGB images, RGB-D images, depth images, point clouds, or meshes) the goal is to predict articulation parameters: *e.g.* articulation type (prismatic *vs.* hinge), segmentation of parts that independently articulate, axis of rotation / translation, points of interaction. Researchers have a) investigated the use of different input modalities [18], [27], [31], [40], b) built datasets for training models [32], c) designed unified output parameterizations [18], and d) designed novel neural architectures and representation [27], [52]. Researchers have also studied directly predicting sites for interaction [31] and trajectories that the robot end-effector should follow [49] to articulate the object. Our work is complementary, and focuses on converting articulation geometry, possibly predicted from any of these past models, into motion plans.

Simulators for studying object articulation have been challenging to build. Most past efforts use manually created *synthetic* scenes: AI2-THOR [23], Sapeins [51], ManipulaTHOR [9], ThreeDWorld [11]. Habitat 2.0 [45] and iGibson [41] improve realism by manually aligning 3D models to real scenes, but are small in size (92 objects in 1 home and 500 objects in 15 homes respectively). Our proposed ArtObjSim simulator is unique in its focus, studying prediction of motion plans for everyday articulated objects, and scale, having 3758 articulated objects spread across 97 unique real world scenes. To our knowledge, ArtObjSim is the largest dataset, to date, for the study of motion planning performance for articulating everyday objects in everyday scenes.

End-to-end RL approaches can also be used to leverage prior experience for fast execution under partial information at test time [14], [26], [50]. However, the large sample complexity of learning policies through RL and the small number of environments available for training has prevented past works to show generalization results in novel environments. By leveraging classical components and scaling up learning, we are able to learn models that generalize to novel objects. **Learning for motion planning** has been used to reduce the runtime of motion planning algorithms: [15], [16], [43]. Strudel *et al.* [43] learn obstacles representations for motion planning, while Ichter *et al.* [15], [16] use learning to bias sampling of states for motion planners. Our use of learning is similarly motivated, but we learn to predict low-dimensional strategies (that can be decoded into full motion plans) for

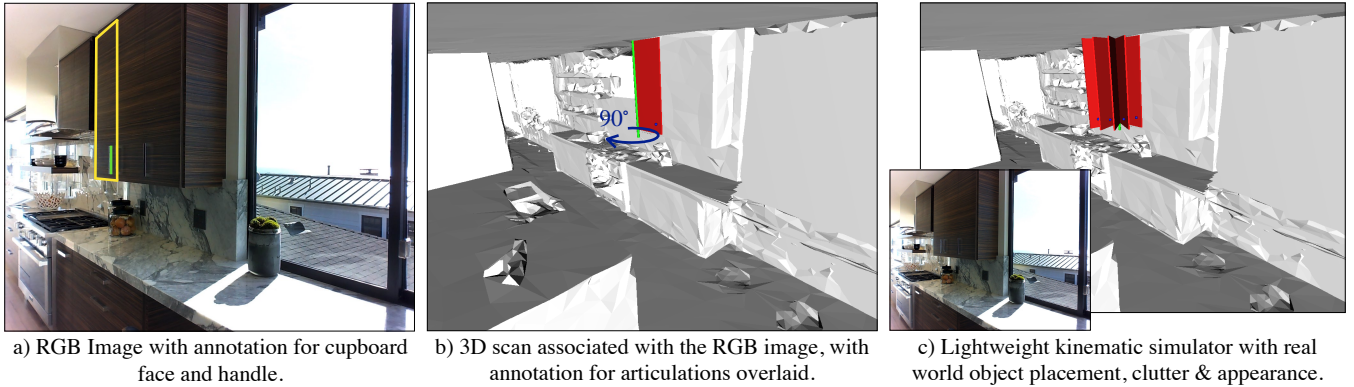


Fig. 2: ArtObjSim development. (a) We annotate RGB images inside 3D scans with 2D articulation geometry. (b) This is lifted to 3D using the underlying 3D geometry. (c) As a result we get simulators that can simulate articulated objects in realistic scenes. See Section III for more details.

TABLE I: Statistics for objects and scenes in ArtObjSim, our proposed simulator for everyday articulated objects in real scenes.

	Train	Val	Test	Total
# Scenes	65	17	15	97
# Unique Object Instances	2538	590	630	3758
# Prismatic (<i>e.g.</i> Drawer)	865	191	211	1267
# Vertical Hinge (<i>e.g.</i> Cabinet)	1325	335	332	1992
# Horizontal Down-hinge (<i>e.g.</i> Oven)	146	40	40	226
# Horizontal Up-hinge (<i>e.g.</i> Toilet lid)	202	24	47	273

constrained motion planning problems from visual input.

III. ARTOBJSIM: A SIMULATOR FOR EVERYDAY ARTICULATED OBJECTS IN REAL SCENES

We introduce ArtObjSim, a lightweight kinematic simulator for articulated objects placed in real scenes. ArtObjSim is built upon the HM3D dataset [37]. HM3D consists of 3D scans of real world environments. It offers both, realistic image renderings from real scenes, and access to the underlying 3D scene geometry. ArtObjSim is made possible through 2D annotations of articulation geometry on images, which are then lifted to 3D to allow for a kinematic simulation of the articulated objects. To our knowledge, ArtObjSim is the first simulator that enables a systematic large-scale study of articulation of everyday objects in real world environments. We describe the steps involved in the construction of ArtObjSim.

Annotating Articulation Geometry on Images. The first step is to annotate 2D articulation geometry on images. 2D articulation geometry includes marking the extent, axis of articulation, articulation type, and interaction locations (handles). We collect annotations in two phases.

In the first phase, we manually walk through the HM3D scenes to find kitchens and bathrooms, and identify locations that show articulated objects. We render out images from different viewpoints from these locations for labeling.

In the second phase, we use an annotation service to obtain the necessary 2D labels. We obtain annotations for the segmentation mask for the front face, handle locations, and articulation type (prismatic *vs.* left hinge *vs.* right hinge *top* hinge *vs.* bottom hinge). See Figure 2(a) for an example annotation. For most rectangular objects (*e.g.* drawers, cupboards, refrigerators) these three together with the underlying

3D information from the mesh are sufficient to deduce the axis of articulation. This doesn’t work for toilets and we get additional labels for the axis of rotations (location where the lid is attached). Toilet lids also don’t have handles, we annotate and use the lid tip as the interaction point.

We manually verify the annotation quality after each phase and fix or reject bad annotations. The annotation procedure is fast and cost effective (average \$0.5 per object instance).

Extracting 3D Articulation Geometry from 2D Annotations. We use the collected 2D annotations, combined with the 3D scene geometry, to obtain a 3D simulation for each articulated object. For each object, we fit a plane to the points within the segmentation mask on the depth image corresponding to the RGB image. This gives us a 3D representation (a 3D rectangle) for the object face that will undergo articulation. We project the 2D handle location onto this 3D plane to obtain the 3D handle location. Articulation parameters are obtained from this 3D representation. We assume that the prismatic objects pull out perpendicular to the face, and the hinged objects rotate about the corresponding edge (top, bottom, left or right) of the 3D rectangle. As noted, toilet lids can’t be approximated as rectangles. We project the annotated 2D axis to the 3D plane. All annotations are converted into the mesh coordinate frame using the transformations for the camera used to render out the image. This defines all that we need to simulate the articulating object in 3D, see Figure 2(c).

ArtObjSim Simulator. As a result of the above two steps, we obtain kinematic simulations for thousands of unique object instances placed in real 3D scenes. Not only can we can simulate the object (*i.e.* how the collision geometry will change as the object articulates or how will the end-effector need to move), we also have a sense of the surrounding 3D geometry of the scene (*i.e.* the counter below the cabinet), and can render out the RGB appearance of the object from multiple different views.

Table I shows dataset statistics. The dataset is diverse with 3758 object instances from across 97 scenes across 10 object categories and 4 articulation types. The dataset also includes a large geometric variety *e.g.* cabinets high up above the counter and oven drawers very close to the ground. This

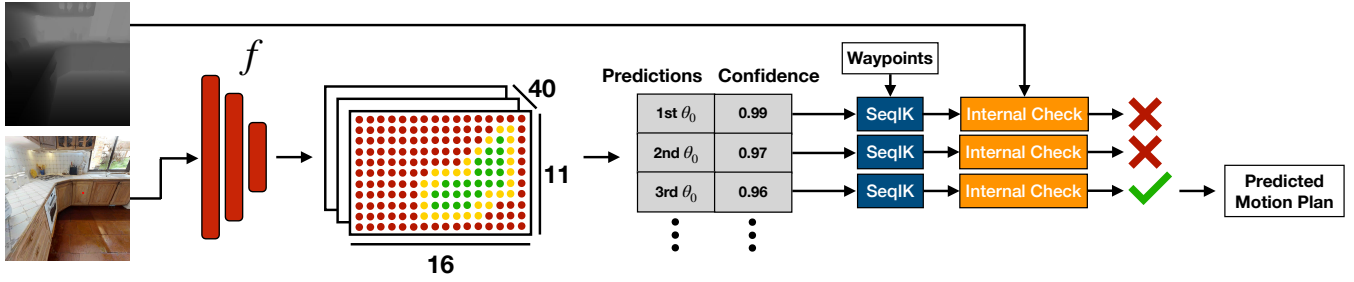


Fig. 3: Overview of MPAO (Motion Plans to Articulate Objects). Given an RGB-D image of the object to be articulated (denoted with a red marker), we use a CNN to predict good initializations for Sequential Inverse Kinematics (SeqIK). SeqIK uses end-effector trajectories to generate motion plans corresponding to each returned high-scoring initializations. Generated plans are tested for deviations from the intended trajectory, and collisions using the depth image. The first plan that succeeds these internal checks is returned.

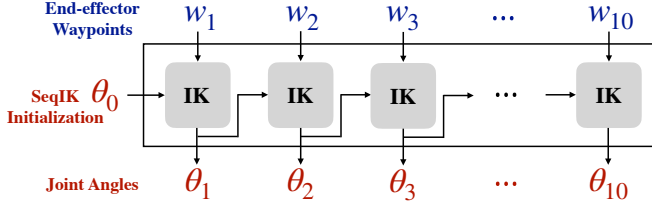


Fig. 4: Sequential Inverse Kinematics (SeqIK). Given an initial joint configuration (θ_0), and a sequence of end-effector pose waypoints, SeqIK uses inverse kinematics (IK) to generate configurations that achieve the given end-effector waypoints. IK for subsequent steps is warm-started with IK solutions from the previous time step.

diversity, along with the fact that these objects are immersed in real scenes makes up problem instances which have not been tackled extensively in the literature.

In Section IV, we will use ArtObjSim to design, train, and evaluate models for predicting motion plans for articulating everyday objects. However, we anticipate ArtObjSim will be useful for many other tasks. For example, predicting articulation parameters or end-effector waypoints from RGB images, or for mining statistics about placement of articulated objects in kitchens to build generative models for scene layout, or for building policies for mobile manipulation.

IV. REPRESENTING AND PREDICTING MOTION PLANS

Given a single RGB-D image pair $[I, D]$ of an articulated object, and a sequence of end-effector poses necessary to articulate the object $[\dots, w_t, \dots]$, our next goal is to predict a motion plan, *i.e.* sequence of joint angles $[\dots, \theta_t, \dots]$ that bring the end-effector in the necessary pose to conduct the desired articulation. Rather than re-solving each new problem instance from scratch using motion planning under partial information, we pursue a machine learning approach that leverages past experience to directly predict motion plans. A straight-forward application of machine learning doesn't work as the predicted plans need to satisfy tight task constraints. Instead, we use machine learning to predict a *strategy* which is decoded into a complete motion plan that adheres to the task constraints at hand. We first describe what strategies are and how they are decoded in Section IV-A and then describe how we use them to predict motion plans from RGB images in Section IV-B.

A. Representing and Decoding Motion Plans

We represent motion plans as the initialization of a deterministic gradient-based solver that optimizes joint angles to get the end-effector in the desired pose.

Our motion plan representation builds upon numerical inverse kinematics methods [28]. Inverse kinematics (IK) is the process of obtaining joint angles that get the end-effector to a given desired pose. Starting from some initial joint angles, a numerical IK solver iteratively updates the joint angles using the Jacobian of the forward kinematics till a solution is found. As we are interested in not one but a sequence of joint angles that track the given end-effector trajectory, we *sequentially* solve a sequence of inverse kinematic problems by initializing the inverse kinematic solver for the t^{th} time-step with the solution from the $(t-1)^{\text{th}}$ time-step. We call this process, Sequential Inverse Kinematics or SeqIK, and show a block diagram in Figure 4.

Thus, SeqIK can be viewed as a *deterministic* process that converts a sequence of end-effector waypoints and an initial joint configuration θ_0 into a sequence of joint angles that realize the end-effector poses. θ_0 can be thought of as *knob* that controls the motion plans that SeqIK generates. Varying θ_0 varies the motion plan generated. We use (SeqIK, θ_0) as our representation for *strategies* that generate motion plans. Our experiments demonstrate that it is a flexible and efficient way to generate motion plans for articulating everyday objects, and outperforms both unconstrained and constrained motion planning approaches.

Note that $\text{SeqIK}+\theta_0$, shorthand for (SeqIK, θ_0) , may not generate feasible motion plans for all inputs θ_0 . Initializing from some θ_0 may not get the end-effector to where we want it to be, others might cause the end-effector to deviate too much from the desired trajectory when interpolating between waypoints, yet others might cause collisions with self or with the environment. We address this issue by *predicting* good θ_0 s from the RGB image showing the articulated object as we describe in Section IV-B.

B. Predicting Motion Plans from Images

Our next step is to predict good initializations θ_0 s for $\text{SeqIK}+\theta_0$ from RGB images. As there can be more than one good θ_0 for each image, we adopt a classification

TABLE II: Motion planning for articulating objects under full information. We measure the success rate and quality of successful plans generated by the different motion planning methods we considered. We note that SeqIK+ θ_0 is able to successfully generate plans quickly. Motion planning, both unconstrained and constrained, obtained a 0% success rate, and hence are omitted from the table, see Section V-A for more details.

Articulation Type	Performance			Speed	
	Success Rate (%)	Translational Deviation (m)	Rotational Deviation (rad)	Median number of initializations	Time (s)
Prismatic (e.g. Drawers)	99.1	0.0005	0.0006	38	9.80
Vertical Hinge (e.g. Cabinets)	63.3	0.0013	0.0026	1118	434.12
Horizontal Down-Hinge (e.g. Dishwasher)	71.8	0.0029	0.0024	896	517.03
Horizontal Up-Hinge (e.g. Toilet lid)	44.7	0.0025	0.0028	1221	13 154.12

approach. We work with a set of initializations Θ . We train a function $f(I, \theta_0)$ that classifies whether or not the use of θ_0 serves as a good initialization for SeqIK to achieve end-effector waypoints $[\dots, w_t, \dots]$ without collisions. We provide details about the initialization set Θ , function f , training data, and loss function to train f .

Initialization set Θ comes from the Cartesian product of a set of robot base positions in \mathbb{R}^3 and a set of 10 arm configurations. We use 704 base positions (sampled in a uniform $1\text{ m} \times 1.5\text{ m}$ 2D grid of base positions at a 10 cm resolution at 4 different heights) and 10 arm configurations, resulting in Θ having 7040 elements. The 10 arm configurations are selected in a data-driven manner, we sample 20 random configurations that satisfy the joint limits, and then select the 10 which when used with SeqIK give us the most successes across the training set.

Function f is realized through a CNN with an ImageNet pre-trained ResNet-34 backbone [13]. We add 2 fully connected layers on the conv5 output to produce a 7040 dimensional representation. This is reshaped into an 80-dimensional spatial output of size 11×16 . This is processed through another 3 convolutional layers to produce a $(10 \cdot 4) \times 11 \times 16$ logits tensor containing 11×16 spatial output for each of the 10 arm configurations at each of the 4 heights.

Training labels are generated by decoding each candidate θ_0 into motion plans using SeqIK, and testing them for end-effector pose deviation, self-collision, collision with the static environment, and collision with the articulating object in our simulator from Section III. Note that while testing the decoded motion plans, we interpolate between consecutive states to simulate how the plan will be executed in practice. This process generates a binary success label for each of the 7040 candidates in Θ . This is used to supervise the logits predicted by f via a binary cross-entropy loss.

Training details. Each articulated object instance in ArtObjSim comes with waypoints and ground truth labels as described above. We render multiple views for each articulated object to generate 40K images to train the function f .

Our full method, *Motion Plans to Articulate Objects (MPAO)*, uses the learned function f to rank candidate initialization in Θ . We go down the ranked list, decode them into motion plans using SeqIK, and return the first *feasible* plan (feasible meaning: accurately tracks the given waypoints and also doesn't collide with self or with the geometry visible in the depth image). See Figure 3 for an overview.

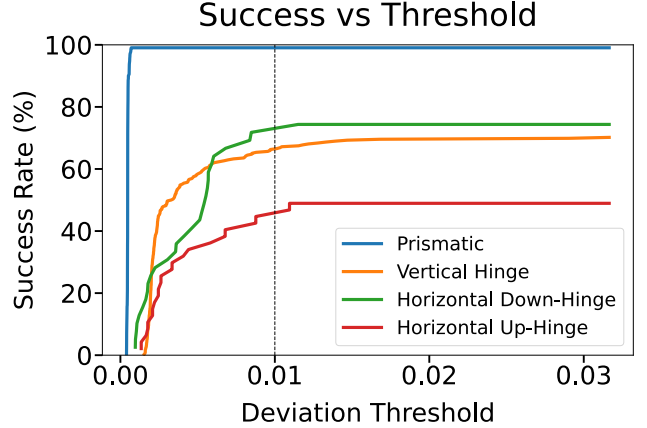


Fig. 5: Success rate as a function of deviation for motion plans generated by our proposed SeqIK. X-axis shows maximum translational and rotational deviation (in m and rad respectively), Y-axis shows the success rate for motion plans generated by SeqIK. SeqIK generates motion plans with low deviation from desired end-effector trajectories.

V. EXPERIMENTS

Our experiments evaluate two aspects: a) the flexibility and decoding efficiency of our proposed motion plan representation from Section IV, and b) how effectively can we leverage RGB images to quickly convert end-effector poses to motion plans. For the former, we make comparisons to motion planning, and for the latter we compare against variations that don't use the RGB image. We also evaluate how our method works with predicted end-effector waypoints.

Experimental Setup. We leverage the geometry and appearance of articulated objects in real scenes in our proposed ArtObjSim simulator for evaluation. We adopt the train, val, and test splits as noted in Table I. All instances from the same scene are in the same set. This allows us to measure how well our models perform on novel held-out object instances. We work with the 7DOF Franka Emika Panda robot. We assume that it can take one of 4 discrete heights (0.25 m, 0.5 m, 1.0 m and 1.5 m). While we reason about where the base should be to conduct the motion, we assume that the base remains fixed during execution. Leveraging base motion to better articulate objects is left to future work.

A. Motion Plan Representation

We evaluate the flexibility and decoding efficiency of our proposed motion planning representation. More specifically, given a 10 time-step end-effector trajectory and complete collision geometry of the situation, this evaluation measures

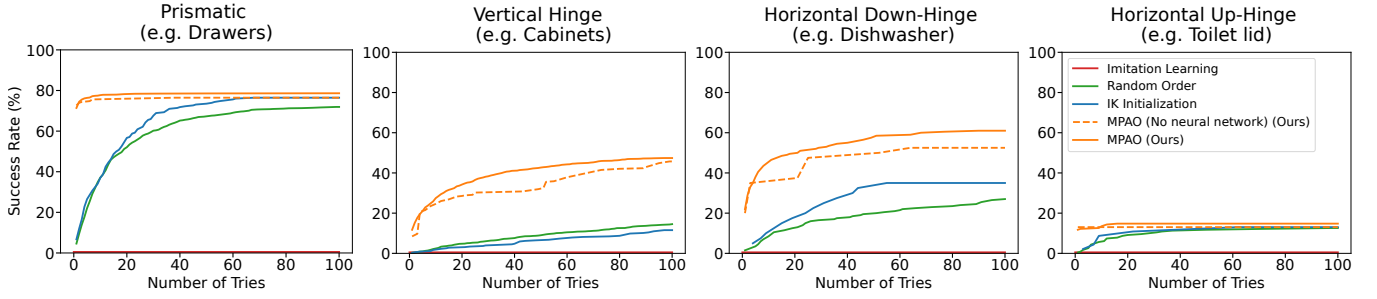


Fig. 6: Motion plan prediction success rate and speed. We show success rate as a function of number of tries for the different articulation types. Our method MPAO, achieves a higher success rate and generates solutions faster than pure search or pure learning methods. We use 0.01 m translational and 0.01 rad rotational tolerance on the end-effector pose to determine success.

the quality of the joint angle trajectory produced by our method. We search for a good initialization $\theta_0 \in \Theta$ for SeqIK and spits out the first solution that doesn’t have collisions (to self, surrounding environment, or the articulating object) and conforms to the given tolerance in end-effector pose.

Metric. A predicted trajectory is considered successful if: a) it conveys the end-effector to the goal pose within 0.01 m and 0.01 rad, b) the resulting end-effector trajectory violates the task constraint by less than 0.01 m in translation and 0.01 rad in rotation for each time step, and c) it doesn’t cause collisions with self, the static environment, or the object as it articulates. Before measuring deviations and collision-checking, we linearly interpolate the joint angle trajectories to bring all joint angle changes to ≤ 0.01 rad.

Results. We report the success rate and time taken by our method for different articulation types in Table II. Prismatic drawers are easy: we can find solutions for 99.1% of the instances to within 0.0005 m and 0.0006 rad deviation, in as little as 10 s of compute while only needing to try a median of 38 initializations. Vertical hinged and horizontal down-hinged objects are harder: we are only able to solve 63.3% and 71.8% instances respectively while also needing to sample many more initializations, taking around 500 s. Toilets are by far the hardest because of the tight space in bathrooms, and also because of the increased collision checking cost due to the non-cuboidal toilet lid geometry. Figure 5 plots the success rate as a function of the deviation (max translational or rotational deviation in m or rad). While we reported success rate at deviation threshold of 0.01 in Table II, Figure 5 shows that most plans returned by SeqIK are much more accurate.

Comparison with other methods. We also compared SeqIK+ θ_0 to two other class of methods: unconstrained and constrained motion planning, neither of which were able to find any successful solutions in a tractable amount of time. For **unconstrained motion planning**, we used RRT-connect [24] to find a path between a start and end joint configuration obtained using inverse kinematics. While this always found a path, without any constraint on the intervening end-effector poses, the path would always violate the 1-DOF constraint imposed by articulated object. This is not surprising as the two poses are quite far from one another. To our surprise, even when these poses are brought

close to one another, by sampling 10 way-points along the trajectory, unconstrained motion planning would still only return solutions that would wildly swing the end-effector around. For **constrained motion planning**, we used the projected state space method from the OMPL library [20], [21], [44]. It would find motion plans that conformed to the task constraint to some extent. However, the minimum translation deviation was 0.02 m, much more than the tolerance level needed in our tasks, resulting again in a 0% success rate. We experimented with many hyper-parameter settings. Some worked better than others, but none were able to return any plans with less than 0.02 m translation deviation.

In summary, SeqIK+ θ_0 is effective at producing joint angles that conform to a given end-effector trajectory. Finding a solution is still computationally expensive as it requires testing many initializations. We address this using the prediction network f . We evaluate it next.

B. Motion Plan Prediction with Known Waypoints

Our next evaluation seeks to measure how quickly and accurately, we can predict motion plans for articulated objects places in novel contexts as observed through RGB-D images. More specifically, given an RGB-D image along with an end-effector trajectory, we measure the success rate of predicting motion plans as a function of planning time. As in Section V-A, we call a predicted motion plan successful if it reaches the goal while violating the task constraint by less than 0.01 m, 0.01 rad and not colliding with self, the environment, and the articulating object. While the metric is the same, the focus of this evaluation is to assess how well methods can cope with partial information from RGB-D observations and their speed of generating solutions.

Comparisons. We compare against other search schemes for finding good θ_0 for SeqIK. These baseline schemes employ the same overall structure as our method (SeqIK decoding followed by filtering based on feasibility), but don’t use any past experience (learned model) to rank initializations. We consider two baselines, *Random Order* and *IK initialization*, and a variant of our method, *MPAO (No neural network)*. We also compare to a pure machine learning approach (*Imitation Learning*) that uses imitation learning to directly predict motion plans. We describe these in more detail below:

- *Random Order* tests whether our learned function f has

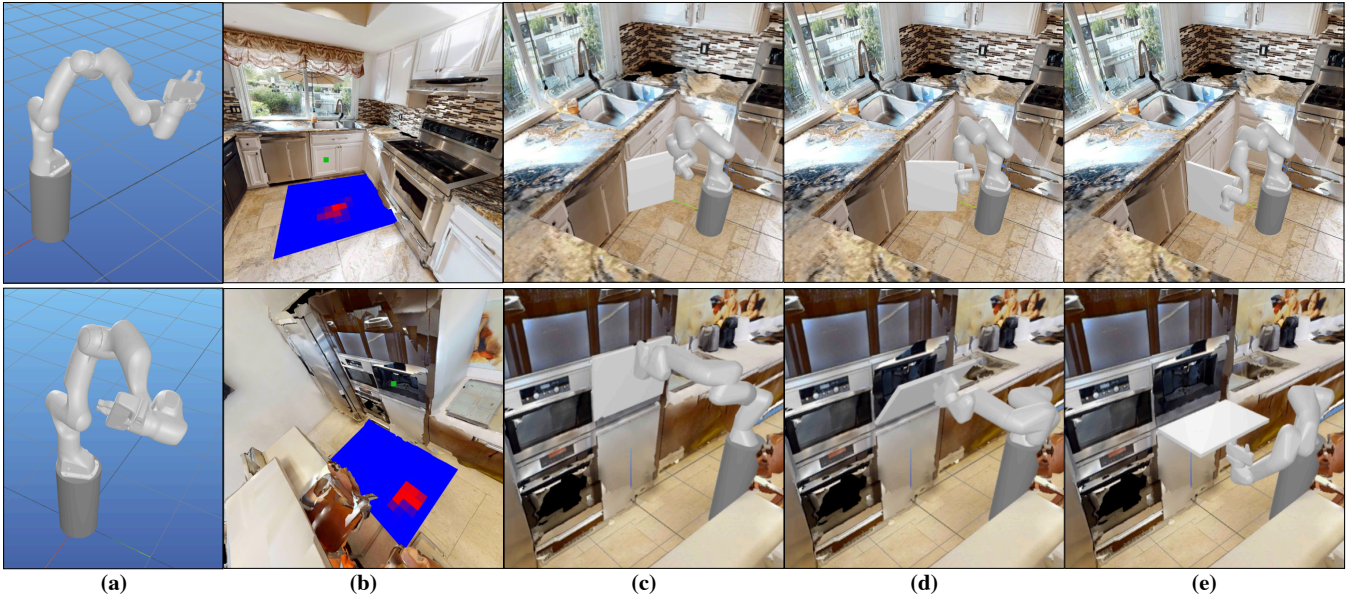


Fig. 7: (a) One of the ten arm joint configurations from Θ used for initialization. (b) Example of an object from the dataset (indicated by the green marker), along with predictions for the configuration shown in (a) overlaid onto the image (warmer colors mean higher score). (c, d, e) Visualizations of a successful execution from one of the high-scoring locations.

learned anything meaningful about which initializations are good, and which initializations are bad. Rather than using our learned function f to sort the set of initializations Θ for each object, we try them in a random order.

- *IK initialization* uses IK to find not only the joint angles but also the base location for the first waypoint. This way, the base location is not restricted to our discrete grid of base positions. After this point, SeqIK is used to obtain a trajectory with a fixed base position, just as for our method.

- *MPAO (No neural network) (Ours)* is a simplification of MPAO that doesn't use the neural network f , but instead ranks initializations in Θ by their success rate on the training set. Though this doesn't use the neural network, it is still *data-driven* in that it leverages experience with past constrained motion planning problems to output plans.

- *Imitation Learning.* We also experiment with a pure imitation learning approach that directly predicts the entire motion plan. The model takes in as input an RGB image (with a dot specifying which articulated object to interact with), as well as 3D waypoints in the camera coordinate frame. As there may be multiple correct ground truth trajectories associated with a given input, the model outputs 100 trajectories, where each trajectory is represented by ten joint configurations (one for each of the waypoints). We use the Hungarian Matching algorithm [25] to assign each predicted trajectory the closest unique ground truth trajectory from the set of ground truth trajectories (as also used for object detector training [7]). We use an L_2 loss on the inferred pairings to train the model. We also jointly train a loss prediction network in order to sort the predictions. At inference time, we employ the same internal check (as used for all other methods) on the ranked predictions before outputting a solution.

Results. Figure 6 presents the success rate for different

methods as a function of total number of solutions tried for novel object instances in the test set. Across all articulation types, our method dominates pure search baselines in success rate and speed. For all categories, we are able to match baseline performance with $10\times$ fewer tries, and obtain more than 25% absolute improvement in success rate for vertical and horizontal down hinges. Even just ranking initializations based on their performance on the training set, *i.e.* MPAO (No neural network), works quite well. This suggests the utility of leveraging past experience for this problem. Our full method, MPAO, boosts performance further and is able to effectively leverage the RGB observation to improve the ranking among solutions.

Imitation learning struggles for this task and yields a 0% success rate on our benchmark. We investigated this further. While imitation learning produced the correct general motion, the motion wasn't precise enough. When evaluated under looser criterion (max deviation of 0.10 m and 0.10 rad vs. 0.01 m and 0.01 rad used in our benchmark), the imitation learning baseline obtains a non-trivial success rate (73.5%, 2.9%, 15.6% and 13.3% respectively for the 4 categories). MPAO comfortably outperforms this even when evaluated under the tighter 0.01 m 0.01 rad criterion (success rates of 78.6%, 47.4%, 61.0% and 14.8% respectively as seen in Figure 6). We also experimented with using the depth image in addition to the RGB image for the imitation learning model, but this did not increase performance.

These experiments together establish the effectiveness of our method at predicting good motion plans. Figure 7 shows an example visualization of motion plans output by MPAO.

C. Motion Plan Prediction with Unknown Waypoints

As a proof-of-concept, we have also integrated MPAO into an overall pipeline that doesn't require known waypoints. We

experimented with drawers. We adapt Mask RCNN [12] to detect and predict drawer faces (segmentation mask) and handle locations (keypoints) using annotations from ArtObjSim. We convert them into end-effector waypoints using the depth image. This by itself gave a median error of 1.6 cm. When using MPAO to track these predicted waypoints, we are able to predict plans that solve 39% drawers to within 0.01 m translation error and 70% to within 0.05 m translation error.

VI. CONCLUSION

We pursued a learning approach that uses past experience to quickly predict motion plans for articulating objects. We collected a large dataset to build ArtObjSim, a simulator that enables a kinematic simulation of everyday objects placed in real scenes. We designed SeqIK+ θ_0 , a fast and flexible way to represent motion plans under end-effector constraints, and trained neural network models that leverage SeqIK+ θ_0 to quickly predict plans for articulating novel objects.

ACKNOWLEDGEMENTS

This material is based upon work supported by DARPA (Machine Common Sense program), an NSF CAREER Award (IIS-2143873), the Andrew T. Yang Research and Entrepreneurship Award, an Amazon Research Award, an NVIDIA Academic Hardware Grant, and the NCSA Delta System (supported by NSF OCI 2005572 and the State of Illinois). We thank Matthew Chang and Aditya Prakash for helpful feedback on experiment design and writing.

REFERENCES

- [1] Ben Abbatematteo, Stefanie Tellex, and George Konidaris. Learning to generalize kinematic models to novel objects. In *CoRL*, pages 1289–1299, 2019.
- [2] JKT Alexander Andreopoulos and John K Tsotsos. A framework for door localization and door opening using a robotic wheelchair for people living with mobility impairments. In *Robotics: Science and systems, Workshop: Robot manipulation: Sensing and adapting to the real world, Atlanta*, 2007.
- [3] Dragomir Anguelov, Daphne Koller, Evan Parker, and Sebastian Thrun. Detecting and modeling doors with mobile robots. In *IEEE International Conference on Robotics and Automation, 2004. Proceedings. ICRA'04. 2004*, volume 4, pages 3777–3784. IEEE, 2004.
- [4] Dmitry Berenson. *Obeying Constraints During Motion Planning*, pages 1–32. Springer Netherlands, 2018.
- [5] Dmitry Berenson, Siddhartha Srinivasa, and James Kuffner. Task space regions: A framework for pose-constrained manipulation planning. *IJRR*, 30(12):1435–1460, 2011.
- [6] Felix Burget, Armin Hornung, and Maren Bennewitz. Whole-body motion planning for manipulation of articulated objects. In *ICRA*, pages 1656–1662, 2013.
- [7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.
- [8] Sachin Chitta, Benjamin Cohen, and Maxim Likhachev. Planning for autonomous door opening with a mobile manipulator. In *2010 IEEE International Conference on Robotics and Automation*, pages 1799–1806. IEEE, 2010.
- [9] Kiana Ehsani, Winson Han, Alvaro Herrasti, Eli VanderBilt, Luca Weihs, Eric Kolve, Aniruddha Kembhavi, and Roozbeh Mottaghi. Manipulator: A framework for visual object manipulation. In *CVPR*, pages 4497–4506, 2021.
- [10] Farbod Farshidian, Edo Jelavic, Asutosh Satapathy, Markus Gifthalder, and Jonas Buchli. Real-time motion planning of legged robots: A model predictive control approach. In *ICHR*, pages 577–584, 2017.
- [11] Chuang Gan, Jeremy Schwartz, Seth Alter, Martin Schrimpf, James Traer, Julian De Freitas, Jonas Kubilius, Abhishek Bhandwaldar, Nick Haber, Megumi Sano, et al. Threedworld: A platform for interactive multi-modal physical simulation. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.
- [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages 2961–2969, 2017.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [14] Daniel Honerkamp, Tim Welschhold, and Abhinav Valada. Learning kinematic feasibility for mobile manipulation through deep reinforcement learning. *IEEE RA-L*, pages 6289–6296, 2021.
- [15] Brian Ichter, James Harrison, and Marco Pavone. Learning sampling distributions for robot motion planning. In *ICRA*, pages 7087–7094. IEEE, 2018.
- [16] Brian Ichter, Edward Schmerling, Tsang-Wei Edward Lee, and Aleksandra Faust. Learned critical probabilistic roadmaps for robotic motion planning. In *ICRA*, pages 9535–9541. IEEE, 2020.
- [17] Advait Jain and Charles C Kemp. Behaviors for robust door opening and doorway traversal with a force-sensing mobile manipulator. Georgia Institute of Technology, 2008.
- [18] Ajinkya Jain, Rudolf Lioutikov, Caleb Chuck, and Scott Niekum. Screwnet: Category-independent articulation model estimation from depth images using screw theory. In *ICRA*, pages 13670–13677, 2021.
- [19] Yiannis Karayiannidis, Christian Smith, Francisco Eli Vina Barrientos, Petter Ögren, and Danica Kragic. An adaptive control approach for opening doors and drawers under uncertainties. *IEEE Transactions on Robotics*, 32(1):161–175, 2016.
- [20] Zachary Kingston, Mark Moll, and Lydia E Kavraki. Sampling-based methods for motion planning with constraints. *Annual review of control, robotics, and autonomous systems*, 1:159–185, 2018.
- [21] Zachary Kingston, Mark Moll, and Lydia E. Kavraki. Exploring implicit spaces for constrained sampling-based planning. *IJRR*, 38(10–11):1151–1178, Sept. 2019.
- [22] Ellen Klingbeil, Ashutosh Saxena, and Andrew Y Ng. Learning to open new doors. In *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2751–2757. IEEE, 2010.
- [23] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. AI2-THOR: An Interactive 3D Environment for Visual AI. *arXiv*, 2017.
- [24] James J Kuffner and Steven M LaValle. RRT-connect: An efficient approach to single-query path planning. In *ICRA*, 2000.
- [25] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- [26] Chengshu Li, Fei Xia, Roberto Martín-Martín, and Silvio Savarese. Hrl4in: Hierarchical reinforcement learning for interactive navigation with mobile manipulators. In *CoRL*, pages 603–616, 2020.
- [27] Xiaolong Li, He Wang, Li Yi, Leonidas J Guibas, A Lynn Abbott, and Shuran Song. Category-level articulated object pose estimation. In *CVPR*, pages 3706–3715, 2020.
- [28] Kevin M Lynch and Frank C Park. *Modern robotics*. Cambridge University Press, 2017.
- [29] Wim Meeussen, Melonee Wise, Stuart Glaser, Sachin Chitta, Conor McGann, Patrick Mihelich, Eitan Marder-Eppstein, Marius Muja, Victor Eruhimov, Tully Foote, John Hsu, Radu Bogdan Rusu, Bhaskara Marthi, Gary Bradski, Kurt Konolige, Brian Gerkey, and Eric Berger. Autonomous door opening and plugging in with a personal robot. In *2010 IEEE International Conference on Robotics and Automation*, pages 729–736. IEEE, 2010.
- [30] Mayank Mittal, David Hoeller, Farbod Farshidian, Marco Hutter, and Animesh Garg. Articulated object interaction in unknown scenes with whole-body mobile manipulation. In *IROS*, 2022.
- [31] Kaichun Mo, Leonidas Guibas, Mustafa Mukadam, Abhinav Gupta, and Shubham Tulsiani. Where2act: From pixels to actions for articulated 3d objects. In *ICCV*, 2021.
- [32] Kaichun Mo, Shilin Zhu, Angel X. Chang, Li Yi, Subarna Tripathi, Leonidas J. Guibas, and Hao Su. PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding. In *CVPR*, 2019.
- [33] Keiji Nagatani and SI Yuta. An experiment on opening-door-behavior by an autonomous mobile robot with a manipulator. In *IROS*, volume 2, pages 45–50. IEEE, 1995.
- [34] Günter Niemeyer and J-JE Slotine. A simple strategy for opening an unknown door. In *Proceedings of International Conference on Robotics and Automation*, volume 2, pages 1448–1453. IEEE, 1997.

- [35] Johannes Pankert and Marco Hutter. Perceptive model predictive control for continuous mobile manipulation. *IEEE RA-L*, pages 6177–6184, 2020.
- [36] L Peterson, David Austin, and Danica Kragic. 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)*, volume 3, pages 2333–2338. IEEE, 2000.
- [37] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexander Clegg, John M Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3D): 1000 large-scale 3d environments for embodied AI. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [38] T. Ruhr, J. Sturm, D. Pangercic, M. Beetz, and D. Cremers. A generalized framework for opening doors and drawers in kitchen environments. In *ICRA*, pages 3852–3858, 2012.
- [39] R. B. Rusu, W. Meeussen, S. Chitta, and M. Beetz. Laser-based perception for door and handle identification. In *International Conference on Advanced Robotics*, pages 1–8, 2009.
- [40] Radu Bogdan Rusu, Wim Meeussen, Sachin Chitta, and Michael Beetz. Laser-based perception for door and handle identification. In *2009 International Conference on Advanced Robotics*, pages 1–8. IEEE, 2009.
- [41] Bokui Shen, Fei Xia, Chengshu Li, Roberto Martín-Martín, Linxi Fan, Guanzhi Wang, Claudia Pérez-D’Arpino, Shyamal Buch, Sanjana Srivastava, Lyne Tchapmi, et al. igibson 1.0: a simulation environment for interactive tasks in large realistic scenes. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7520–7527. IEEE, 2021.
- [42] Jean-Pierre Sleiman, Farbod Farshidian, Maria Vittoria Minniti, and Marco Hutter. A unified mpc framework for whole-body dynamic locomotion and manipulation. *IEEE RA-L*, pages 4688–4695, 2021.
- [43] Robin Strudel, Ricardo Garcia, Justin Carpentier, Jean-Paul Laumond, Ivan Laptev, and Cordelia Schmid. Learning obstacle representations for neural motion planning. *arXiv preprint arXiv:2008.11174*, 2020.
- [44] Ioan A. Şucan, Mark Moll, and Lydia E. Kavraki. The Open Motion Planning Library. *IEEE Robotics & Automation Magazine*, 19(4):72–82, December 2012. <https://ompl.kavrakilab.org>.
- [45] Andrew Szot, Alex Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Chaplot, Oleksandr Maksymets, Aaron Gokaslan, Vladimir Vondrus, Sameer Dharur, Franziska Meier, Wojciech Galuba, Angel Chang, Zsolt Kira, Vladlen Koltun, Jitendra Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0: Training home assistants to rearrange their habitat. In *NeurIPS*, 2021.
- [46] N. Vahrenkamp, T. Asfour, and R. Dillmann. Robot placement based on reachability inversion. In *ICRA*, pages 1970–1975, 2013.
- [47] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J Guibas. Normalized object coordinate space for category-level 6d object pose and size estimation. In *CVPR*, pages 2642–2651, 2019.
- [48] Xiaogang Wang, Bin Zhou, Yahao Shi, Xiaowu Chen, Qingping Zhao, and Kai Xu. Shape2motion: Joint analysis of motion parts and attributes from 3d shapes. In *CVPR*, 2019.
- [49] Ruihai Wu, Yan Zhao, Kaichun Mo, Zizheng Guo, Yian Wang, Tianhao Wu, Qingnan Fan, Xuelin Chen, Leonidas Guibas, and Hao Dong. VAT-Mart: Learning visual action trajectory proposals for manipulating 3d articulated objects. In *ICLR*, 2022.
- [50] Fei Xia, Chengshu Li, Roberto Martín-Martín, Or Litany, Alexander Toshev, and Silvio Savarese. Relmogen: Leveraging motion generation in reinforcement learning for mobile manipulation. *ICRA*, 2021.
- [51] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X. Chang, Leonidas J. Guibas, and Hao Su. Sapien: A simulated part-based interactive environment. In *CVPR*, pages 11097–11107, 2020.
- [52] Vicky Zeng, Timothy E Lee, Jacky Liang, and Oliver Kroemer. Visual identification of articulated object parts. In *IROS*, pages 2443–2450, 2021.