

Bootstrapping Motor Skill Learning with Motion Planning

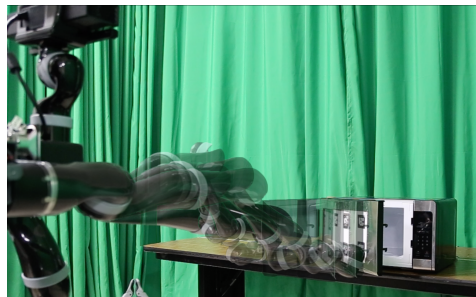
Ben Abbatematteo*, Eric Rosen*, Stefanie Tellex, George Konidaris
Department of Computer Science
Brown University
Providence, RI
{babbatem,er35,stefie10,gdk}@cs.brown.edu

Abstract—Learning a robot motor skill from scratch is impractically slow; so much so that in practice, learning must be bootstrapped using a good skill policy obtained from human demonstration. However, relying on human demonstration necessarily degrades the autonomy of robots that must learn a wide variety of skills over their operational lifetimes. We propose using kinematic motion planning as a completely autonomous, sample efficient way to bootstrap motor skill learning for object manipulation. We demonstrate the use of motion planners to bootstrap motor skills in two complex object manipulation scenarios with different policy representations: opening a drawer with a dynamic movement primitive representation, and closing a microwave door with a deep neural network policy. We also show how our method can bootstrap a motor skill for the challenging dynamic task of learning to hit a ball off a tee, where a kinematic plan based on treating the scene as static is insufficient to solve the task, but sufficient to bootstrap a more dynamic policy. In all three cases, our method is competitive with human-demonstrated initialization, and significantly outperforms starting with a random policy. This approach enables robots to efficiently and autonomously learn motor policies for dynamic tasks without human demonstration.

I. INTRODUCTION

Robots require motor policies for interacting with objects in their environment. For example, a robot butler may need a motor skill that enables it to open a drawer to fetch utensils for a table, for setting each element of the table, and for pouring wine. While it is safe to assume that a robot will have an accurate kinematic model of its own body, it is unlikely to have a dynamics model of every object it will ever encounter. This lack of knowledge means that the robot will have to learn how to manipulate the world around it [20].

Reinforcement learning (RL) provides a framework for robots to acquire motor policies without explicitly modeling the unknown world, but model-free RL methods like policy search [9] have high sample-complexity, and often fail to learn a reasonable policy from random initialization. Supervised approaches for policy learning like Learning From Demonstration (LfD) [2] can encode human prior knowledge by imitating expert examples, but do not support optimization in new environments. Combining RL with LfD is a powerful method for reducing the sample complexity of policy search, and is often used in practice [23, 33, 42, 6]. However,



(a)



(b)

Fig. 1: A robot using our method to autonomously learn to close a microwave. (a) The robot uses a motion planner to generate an initial attempt at closing the microwave door using a kinematic model of the microwave. The resulting plan is unable to fully close the microwave door because of the robot’s limited reach. (b) After bootstrapping a motor skill with the trajectory from (a), the robot learns a motor skill that gives the door a push, exploiting its dynamics to fully close the microwave.

this approach typically requires a human demonstrator for initialization, which fundamentally limits the autonomy, and therefore utility, of a robot which may need to acquire a wide range of motor skills over its operational lifetime. More recently, model-based control techniques (including Model Predictive Control [16, 30] and LQR [23]) have been proposed as exploration methods for policy search; these methods still require human demonstrations or complete dynamic models of both the robot and every object in the

*These authors contributed equally to this work

scene.

We propose the use of kinematic motion planning to initialize motor skill policies. While previous work have leveraged sample-based motion planners for learning motor skills [40, 15, 13], they only focus on either free-space motions or do not learn a closed-loop controller. To our knowledge, this is the first use of motion planning to provide initial demonstrations for learning closed-loop motor skill policies by leveraging estimated object kinematics. Motion planning algorithms generate collision-free behavior and generalize to novel scenarios when the robot has a good kinematic model of itself and the object it aims to manipulate, making it useful for tasks like pick-and-place. We show that given a (potentially approximate, and readily estimated) kinematic description of the environment and the robot, off-the-shelf motion planning algorithms can generate feasible (potentially successful but inefficient) initial trajectories (Figure 1a) to bootstrap an object-manipulation policy that can subsequently be optimized using policy search (Figure 1b). This framework enables the robot to automatically produce its own demonstrations for effectively learning and refining object manipulation policies. Our work enables the robot to realize the benefits of an initial demonstration fully automatically using kinematic planning, requiring no human involvement.

To evaluate our method, we used two different motor policy classes (Dynamic Movement Primitives (DMPs) [12] and deep neural networks [24]). We chose these two different motor policy classes because deep neural networks are extremely expressive in what policies they can represent, but are extremely sample inefficient compared to structured motor primitives like DMPs, and we aim to evaluate how our method performs in both contexts. Using these motor policy classes, we compared bootstrapping with motion planning against learning from scratch in three simulated experiments, and against human demonstrations in real hardware experiments. Human demonstrations provide a baseline for how effective these motor policies can do when bootstrapped with high-quality demonstrations, and learning from scratch provides a baseline for how difficult the task is without any prior information about the task. In the first two experiments—opening a drawer with a dynamic movement primitive representation, and closing a microwave door with a deep neural network policy—we show that motion planning using a kinematic model produces a reasonable initial policy, although suboptimal compared to a supervised human demonstration, that learning adapts to generate efficient, dynamic policies that exploit the dynamics of the object being manipulated. We also show how our method can bootstrap a motor skill for the challenging dynamic task of learning to hit a ball off a tee, which involves precise and agile movement. In that case, treating the objects in the scene as static and applying kinematic motion planning succeeds in generating a policy that makes contact with the ball, which is sufficient to bootstrap a more dynamic policy that learns to hit the ball several feet. Our method is competitive with human-demonstrated

initialization, but requires no human demonstration. It serves as a suitable starting point for learning, and significantly outperforms starting with a random policy. This approach enables robots to efficiently and autonomously learn motor policies for dynamic tasks without human demonstration.

In summary, our contributions are:

- 1) A fully-autonomous paradigm for policy search, in which an agent first uses goal-directed kinematic planning to devise feasible solution trajectories for itself and objects in a scene.
- 2) A novel algorithm that autonomously generates a set of initial demonstrations for object manipulation via kinematic planning in object and robot configuration spaces.
- 3) An empirical evaluation of this algorithm in which we employ model-free policy optimization after bootstrapping, demonstrating that our method’s performance is comparable to human expert demonstration and superior to random initialization in hardware and simulation tasks: closing a microwave, opening a drawer, and hitting a ball off a tee.

II. BACKGROUND

Our goal is to efficiently and autonomously learn robot motor skill policies. To do so, we develop an approach that uses kinematic motion planning to generate initial trajectories, fits a policy to those trajectories using behavioral cloning, and subsequently optimizes that policy via policy search. We now briefly describe policy search, policy representations, learning from demonstration, and motion planning.

A. Policy Search

Policy search methods [9] are a family of model-free reinforcement learning algorithms that search within a parametric class of policies to maximize reward. Formally, given a Markov Decision Process $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma \rangle$, the objective of policy search is to maximize the expected return of the policy π_θ :

$$\max_{\theta} \mathbb{E}_{\mathcal{M}, \pi_\theta} \left[\sum_{t=0}^T \gamma^t r_t \right]. \quad (1)$$

These approaches can learn motor skills through interaction, and therefore do not require an explicit environment model, and are typically agnostic to the choice of policy class (though their success often depends on the policy class having the right balance of expressiveness and compactness). However, their model-free nature leads to high sample complexity, which often makes them infeasible to apply directly to robot learning problems.

Policy representation describes the class of functions used as the mapping from states to actions. Our approach is agnostic to policy representation; we demonstrate our approach enables efficient learning using two different common policy representations: Dynamic Movement Primitives and neural network controllers.

Dynamic Movement Primitives [12, 32] are a description of a non-linear second-order differential equation that exhibits attractor dynamics modulated by a learnable forcing function. DMPs are a popular representation for motor policies because they are parameter-efficient, can express both point and limit cycle attractors, enable real-time computation, and exhibit temporal invariance that does not effect the attractor landscape. We refer the reader to the work of Ijspeert et al. [12] for a more formal introduction to DMPs. If we have n joints we wish to control, we can model control for each joint independently with n DMPs for each one. Therefore, the multiple joints of a robot are only coupled through time, which makes this representation very compact.

Neural network controllers have received significant attention in recent years; they are able to learn hierarchical feature representations for approximating functions (in our case, motor skills) operating on high-dimensional input such as robot sensor data. They are more expressive than restricted policy classes such as DMPs and can operate directly on high-dimensional state spaces (e.g. images), yet they typically exhibit higher sample complexity [24].

Learning from Demonstration methods [2, 34] broadly consist of two families of approaches that either mimic (Behavioral Cloning) or generalize (Inverse Reinforcement Learning) the exemplified behavior. Inverse reinforcement learning methods seek to estimate a latent reward signal from a set of demonstrations; we assume a given reward function, and omit a discussion of inverse reinforcement learning methods here.

Behavioral cloning methods [3, 31, 11] attempt to directly learn a policy that reproduces the demonstrated policies. Given a dataset of expert demonstrations D , the objective of behavioral cloning is: $\max_{\theta} \sum_{(s,a) \in D} \pi_{\theta}(a|s)$.

These methods often result in policies with undesirable behavior in states not observed during demonstrations, though this can be addressed with interactive learning [35, 29, 39]. In our approach, the existence of a reward function enables the agent to learn robust behavior in states outside of the initial training distribution. Moreover, our experiments demonstrate our approach’s ability to extrapolate beyond suboptimal initial demonstrations.

Many approaches investigate the incorporation of human-provided demonstrations into policy search to drastically reduce sample complexity via a reasonable initial policy and/or the integration of demonstrations in the learning objective [32, 18, 6, 33, 42, 41, 23].

B. Motion Planning

The pose of an articulated rigid body can be defined by the state of each of its movable joints. The space of these poses is called the configuration space \mathcal{C} [27]. Motion planning is the problem of finding a path (sequence of poses) through configuration space such that the articulated object is moved to a desired goal configuration, without encountering a collision.

While there exist many different families of motion planning algorithms, such as geometric, grid-based, and probabilistic road maps [22], they all operate in a similar fashion: given a configuration space \mathcal{C} and start and goal joint configurations $q_0, q^* \in \mathcal{C}$, return a valid path of joint configurations $\{q_t\}_{t=0}^T$ between the start and end configurations. We focus on sample-based motion planning approaches.

Probabilistic motion planners provide a principled approach for quickly generating collision-free robot trajectories. However, online replanning is expensive, and kinematic motion planners are only as effective as their kinematic models are accurate: they generate trajectories directly, and thus cannot be improved through subsequent interaction and learning. Furthermore, kinematic planners produce trajectories that only account for kinematics, not dynamics: they explicitly do not account for forces involved in motion, such as friction, inertial forces, motor torques, etc, which are important for effectively performing contact-rich, dexterous manipulation.

The process of computing the position and orientation $p \in SE(3)$ of a link in a kinematic chain for a given joint variable setting (a point in configuration space) is termed *forward kinematics*. Inversely, computing a configuration to attain a specific end effector pose p is termed *inverse kinematics*. We denote the forward kinematics functions $p = f(q)$.

III. BOOTSTRAPPING SKILL LEARNING WITH MOTION PLANNING

Our methodology is inspired by how humans generate reasonable first attempts for accomplishing new motor tasks. When a human wants to learn a motor skill, they do not start by flailing their arms around in a random fashion, nor do they require another person to guide their arms through a demonstration. Instead, they make a rough estimate of how they want an object to move and then try to manipulate it to that goal. For example, before being able to drive stick shift, a human must first learn how to manipulate a gear shifter for their car. Just by looking at the gear shifter, humans can decide (1) what they should grab (the shaft), (2) where they want the shaft to go (positioned in a gear location), and (3) how the shaft should roughly move throughout the action (at the intermediate gear positions). Similarly, a robot that has a good kinematic model of itself, and a reasonable kinematic model of the object it wishes to manipulate, should be able to form a motion plan to achieve the effect it wishes to achieve.

That plan may be inadequate in several ways: its kinematic model may be inaccurate, so the plan does not work; object dynamics (like the weight of a door, or the friction of a joint) may matter, and these are not represented in a kinematic model; and a feasible and collision-free kinematic trajectory may not actually have the desired effect when executed on a robot interacting with a real (and possibly novel) object. But such a solution is a *good start*; we therefore propose to use it to bootstrap motor skill learning.

Our approach, outlined in Figure 2, leverages the (partial) knowledge the robot has about its own body and the object it is manipulating to bootstrap motor skills. Our method

first assumes access to the configuration space of the robot, denoted as \mathcal{C}_R , as well as its inverse kinematics function f_R^{-1} . This assumption is aligned with the fact that the robot often has an accurate description of its own links and joints and how they are configured during deployment. However, the world is comprised of objects with degrees of freedom that can only be inferred from sensor data. Therefore, our approach only assumes access to estimated kinematics of the object to be manipulated, in the form of configuration space \mathcal{C}_O and forward kinematics f_O . Recent work has shown that estimating these quantities for novel objects from sensor data in real environments is feasible [1, 25], though state-of-the-art estimates still include noise.

Finally, our approach assumes that the task goal can be defined in terms of kinematic states of the robot and environment. Examples of such tasks include pick-and-place, articulated object manipulation, and many instances of tool use. (Note that this requirement fails to capture reward functions defined in terms of force, for example exerting a specific amount of force in a target location.) Such a goal, together with object and robot kinematic descriptions, enables us to autonomously generate useful initial trajectories for policy search.

Our approach is outlined in Algorithm 1, and can be broken down into five main steps: 1) collect initial trajectory(s) from a motion planner using estimated object kinematics, 2) fit a policy with these initial trajectories, 3) gather rollouts to sample rewards for the current policy based on the kinematic goal, and 4) update the policy parameters based on the actions and rewards, 5) repeat steps 3-4.

Algorithm 1 Planning for Policy Bootstrapping

- 1: **procedure** PPB($C_R, f_R^{-1}, C_O, f_O, q_O^*$)
 - 2: $D \leftarrow \emptyset$
 - 3: **for** 0 to N **do**
 - 4: $D \leftarrow \text{InitialMPDemos}(C_R, f_R^{-1}, C_O, f_O, q_O^*) \cup D$
 - 5: $\theta \leftarrow \text{FitPolicy}(D_0, \dots, D_N)$
 - 6: **for** 0 to E **do**
 - 7: $T_0, \dots, T_n \leftarrow \text{Rollout}(\pi, \theta, q_O^*)$
 - 8: $\theta \leftarrow \text{UpdatePolicy}(T_1, \dots, T_n, \theta)$
-

Algorithm 2 Initial Motion Plan Demos

- 1: **procedure** INITIALMPDEMOS($C_R, f_R^{-1}, C_O, f_O, q_O^*$)
 - 2: $T_O \leftarrow \text{MotionPlanner}(C_O, q_O^*)$
 - 3: $g \leftarrow \text{EstimateGrasp}(C_O, f_O)$
 - 4: $eepath \leftarrow \text{GraspPath}(T_O, C_O, f_O, g)$
 - 5: $T_R \leftarrow \text{MotionPlanner}(C_R, eepath, f_R^{-1})$
 - 6: **return** T_R
-

A. Initial Trajectories from Motion Planner

To fit our policy class, we choose to sample N demonstrations using a motion planner. Our approach for collecting

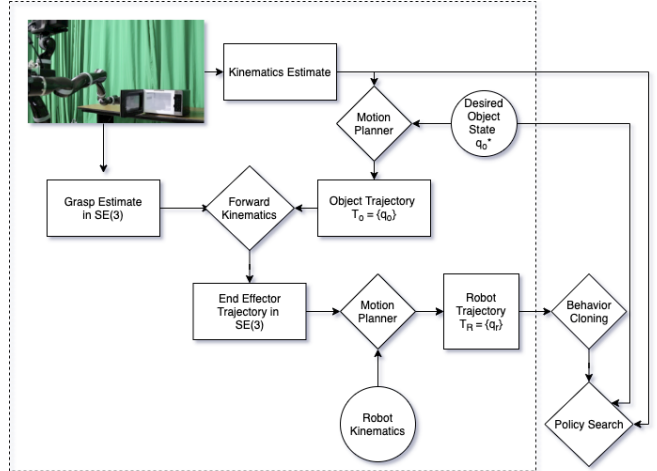


Fig. 2: **System overview** illustrating our proposed framework for generating demonstrations with a motion planner and subsequently performing policy search. The dashed box contains the steps from Algorithm 2.

initial demonstrations with a kinematic motion planner is outlined in Algorithm 2, and proceeds as follows. First, we use a motion planner to find a path through the object’s configuration space \mathcal{C}_O that moves the object from its initial state to a goal state q_O^* using an off-the-shelf motion planner. This produces a joint trajectory in object configuration space, T_O , which transforms the object from its current joint configuration to the desired one.

We then estimate a grasp point on the object to designate the contact point for the robot during manipulation. This can be done by either generating candidate grasps using off-the-shelf grasping algorithms [10, 28] or choosing a part semantically. This produces a local 6D pose, g , that represents where the robot should grasp the object during manipulation.

We then use the grasp point g , object joint trajectory T_O , and the object’s forward kinematics f_O to generate the series of 6D Cartesian poses that the grasp point g will go through as the object proceeds through T_O . This produces a series of 6D Cartesian poses, $eepath$, which the robot end effector must go through, assuming a fixed grasp pose to the object.

Finally, we solve for a path in robot joint space that achieves the end effector path in Cartesian space using off-the-shelf sample-based motion planners, using the robot’s inverse kinematics f_R^{-1} and the sequence of end-effector poses $eepath$.

Note that in our experiments, motion plans were generated offline, rather than recomputed online based on the object’s tracked state, but online motion planning is a trivial extension.

B. Fitting a Policy to a Demonstration

After collecting initial demonstrations from the motion planner, D , we can bootstrap our motor policy by initializing

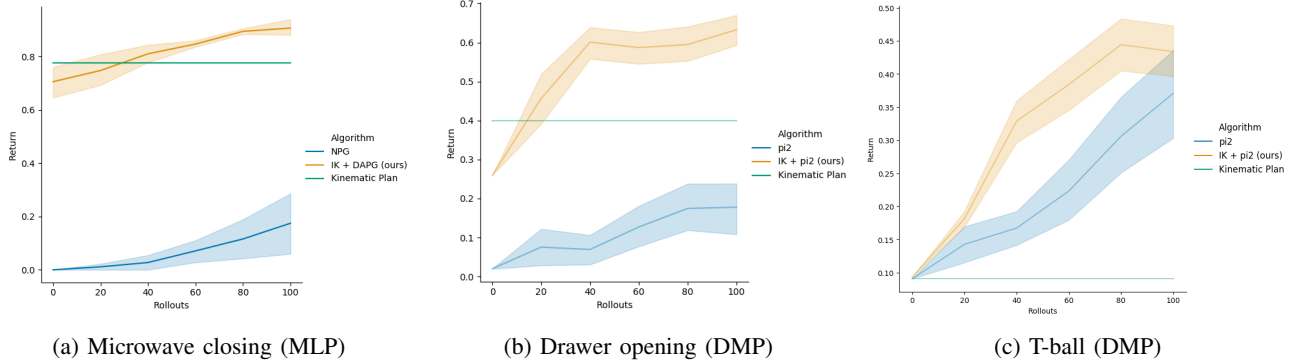


Fig. 3: **Simulation Results.** a) Comparison of our method optimized with DAPG against Natural Policy Gradient starting with a random policy in a microwave closing task using Gaussian multi-layer perceptron policies. b) Comparison of our method against PI²-CMA starting with a random policy in a drawer opening task with DMP policies. c) Our method compared with PI²-CMA with a initially random policy in t-ball with DMP policies. Results are shown as mean and standard error of the normalized returns aggregated across 20 random seeds.

the parameters to the policy θ . We can initialize a parameterized motor policy using any behavioral cloning technique; in practice, for DMPs, we use Locally Weighted Regression [36], and for neural networks, we maximize the likelihood of the demonstration actions under the policy.

C. Policy Search with Kinematic Rewards

To improve the motor policies after bootstrapping, we can perform policy search based on the given (kinematic) reward function. Specifically, we choose a number of epochs E to perform policy search for. For each epoch, we perform an iteration of policy search by executing the policy and collecting rewards based on the goal q_O^* . We define our reward functions using estimated object states q_O and goal states q_O^* , and add a small action penalty.

IV. EXPERIMENTS

The aim of our evaluation was to test the hypothesis that motion planning can be used to initialize policies for learning from demonstration without human input. We tested this hypothesis in simulation, against learning from scratch, and on real hardware, against human demonstrations, on three tasks: microwave-closing, drawer-opening, and t-ball. We note that we do not show asymptotic performance because our emphasis is on learning on real hardware from a practical number of iterations. All the components of the motion planning problem - state sampler, goal sampler, distance metrics, etc. - are reused between problems without modification.

A. Simulation Experiments

The aim of our evaluation was to test the hypothesis that motion planning can be used to initialize policies for learning from demonstration without human input. We tested this hypothesis in simulation against learning from scratch, and on real hardware against human demonstrations, on three tasks: microwave-closing, drawer-opening, and t-ball.

B. Simulation Experiments

We used PyBullet [8] to simulate an environment for our object manipulation experiments. We used URDFs to instantiate a simulated 7DoF KUKA LBR iiwa7 arm and the objects to be manipulated, which gave us ground-truth knowledge of the robot and object kinematics. For all our simulated experiments, we compared implementations of our method against starting with a random policy.

For all three tasks, the state was represented as $s_t = [q_R, q_O]^T$ where q_R denotes robot configuration and q_O denotes object configuration. The action space A was commanded joint velocity for each of the 7 motors. The reward at each timestep r_t was given as:

$$r_t = -c \|q_O^* - q_O\|_2^2 - a_t^T R a_t, \quad (2)$$

where q_O denotes the object state at time t , q_O^* denotes desired object state, and a_t denotes the agent's action. We set $c = 60$ and $R = I \times 0.001$ for all experiments. As such, maximum reward is achieved when the object is in the desired configuration, and the robot is at rest.

Our first simulated task was to close a microwave door, which consisted of three parts: a base, a door, and a handle. The pose of the handle was used for the EstimateGrasp method in Algorithm 2. The robot was placed within reaching distance of the handle when the microwave door was in an open position, but was too far to reach the handle in its closed configuration. Thus, the agent was forced to push the door with enough velocity to close it. We used Gaussian policies represented as multi-layer perceptrons with two hidden layers of sizes (32,32) in this experiment. The randomly initialized policy was optimized with natural policy gradient [17]. Ten demonstrations were generated by perturbing the start state and initial kinematic plan with Gaussian noise. The behavior cloning was performed by maximizing likelihood over the demonstration dataset for 10 epochs. Our pretrained policy was optimized using Demo Augmented Policy Gradient [33],

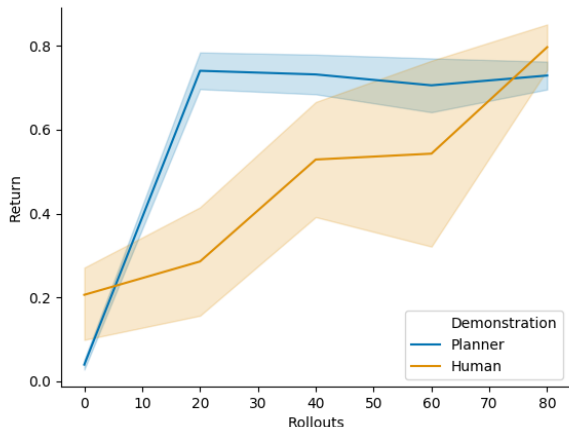


Fig. 4: **Hardware experiment** comparing our initialization scheme with human demonstration. Results are shown as mean and standard error, aggregated across three random seeds.

which essentially adds the behavior cloning loss to the natural policy gradient loss, annealing it over time. This ensures that the agent remains close to the demonstrations early in learning, but is free to optimize reward exclusively as learning progresses. Results are shown in Figure 3a.

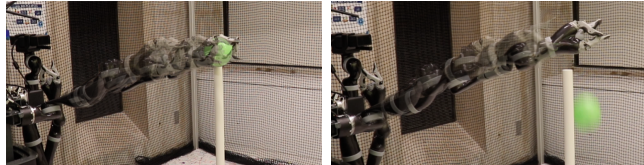
The second simulated task was to open a drawer. This task required the agent to grasp the drawer’s handle and pull the drawer open. Again, the pose of the object’s handle was used for EstimateGrasp method in our algorithm. In this experiment, we used DMP policies. The weights, goals, and speed parameters of the policies were optimized using PI²-CMA [38]. We used 32 basis functions for each of the DMPs. The pretrained policy was initialized using Locally Weighted Regression (LWR) [36] with a single demonstration. The results of this experiment are shown in Figure 3b.

The third simulated task was to hit a ball off a tee. The ball started at rest on top of the tee. The pose of the ball was used in the EstimateGrasp method. The object state was defined as the object’s y position relative to its initial pose. This experiment again used DMP policies initialized with LWR and optimized with PI²-CMA. The results of this experiment are visualized in Figure 3c.

The results of our simulated tasks can be found in Figure 3. Across all three tasks, we observe that policies initialized with our method dramatically outperform starting learning with a random policy. This confirms our hypothesis that using motion planning to generate demonstrations significantly speeds the acquisition of motor skills in challenging tasks like articulated object manipulation and t-ball.

C. Real-world Experiments

For all our real-world experiments, we used a 7DoF Jaco arm [4] to manipulate objects (Figure 1). We used ROS and MoveIt![7] as the interface between the motion planner



(a) Motion Plan Demonstration (b) Bootstrapped from (a)

Fig. 5: **Real-world Ball Hitting** Images comparing an autonomous motion plan demonstration generated from our method vs. a bootstrapped motor skill initialized with that demonstration for a real-world robot hitting a ball off a tee as far as possible. Qualitatively, the bootstrapped motor skill outperforms the initial demonstration by learning to take advantage of the latent tasks dynamics. Videos can be found in our supplemental video. (a) A demonstration provided by the motion planner, which moves linearly towards the ball (b) A motor skill bootstrapped by the motion planner demonstration that learns an agile swooping motion to take advantage of the balls dynamics.

(RRT* [19] in our experiments) and robot hardware. For all real-world experiments, we compared implementations of our method against bootstrapping with a human demonstration, which we supplied. We acknowledge this potential bias in expert trajectories, and qualify our decision by only training on human demonstrations that at least accomplished the task. To collect human demonstrations, we had an expert human teleoperate the robot with joystick control to perform the task. For all tasks, the state space, action space, and reward were defined in the same way as in our simulated results (Section IV-B). Both experiments used DMP policies initialized with LWR [36] and optimized with PI²-CMA [38] with 10 basis functions for each of the DMPs.

Our first real-world task was to close a microwave door, similar to the one described in our simulated domain (Section IV-B). As in the simulated microwave task, we used the pose of the handle for the EstimateGrasp method in Algorithm 2, and also the robot was similarly placed such that it was forced to push the door with enough velocity to close. We placed an AR tag on the front-face of the microwave to track the microwave’s state using a Kinect2. Results are shown in Figure 4. We observe that the human demonstration is better than the one produced by the motion planner, which we credit to the fact that the motion of the door was heavily influenced by the dynamics of the revolute joint which the motion planner did not account for. Nonetheless, both policies converge to a similar final performance, with our method converging slightly faster. Note the importance of the policy search phase: the motion planner alone is insufficient for performing the task efficiently.

Our second real-world task was to hit a ball off a tee as far as possible (Figure 5). Similar to our simulated task, the ball started at rest on top of the tee. The pose of the ball was used in the EstimateGrasp method. The object state was defined as the object’s y position relative to its initial pose.

We placed scotchlite-reflective tape on the surface of the ball and conducted our experiments within an OptiTrack motion-capture cage to track the object pose. We observe that when using a motion planner to hit the ball, it moves the bat in a linear motion to make contact, therefore transferring only horizontal motion to the ball. We qualitatively observe that during policy search, the robot learns a dynamic policy that accounts for the dynamics of the ball by applying force under the ball to “scoop” the ball upwards and forwards.

V. RELATED WORK

To our knowledge, our method is the first to use an object’s estimated kinematics in conjunction with a known robot dynamics model to bootstrap motor policy learning, and we discover and discuss important problems that are only introduced when leveraging policy-learning algorithms, behavioral-cloning, and motion planning algorithms to do so. In this section, we discuss relevant approaches to motor skill learning.

Recently, Model-Predictive Control (MPC) has been used in the context of imitation learning and reinforcement learning to address the high sample complexity of policy search [16, 30]. These approaches require a priori object dynamics, or human demonstrations to fit learned models; in contrast, our approach requires only object kinematics, which are much more readily estimated from visual data at runtime [1, 25]. As such, our approach enables the learning of manipulation skills to be more autonomous than existing MPC-based methods. Tosun et al. [40] proposed a neural network model for generating trajectories from images, using a motion planner during training to enable the robot to generate a trajectory with a single forward pass at runtime. While this approach uses a motion planner for behavior cloning, it stops short of optimization to improve the resulting policy. In contrast, our method uses object kinematics to produce initial trajectories, while Tosun et al. [40] only use the robot’s kinematic model, which is insufficient when the task is to manipulate an object to a specific joint configuration.

While classic robot motor learning papers [3] leverage the known kinodynamics of the robot, they do not discuss kinematics of external objects or grasp candidates to bootstrap motor policies for object manipulation. We emphasize that we cannot form dynamic plans in the problem setting we are interested in: objects with unknown a priori dynamics.

Kurenkov et al. [21] proposed training an initially random RL policy with an ensemble of task-specific, hand-designed heuristics. This improves learning but the initial policy is still random, yielding potentially unsafe behavior on real hardware, and delaying convergence to a satisfying policy. By contrast, we choose to initialize the policy with demonstrations from a kinematic planner, ensuring feasibility, safety, and rapid learning. Moreover, we argue that motion planning is the principled heuristic to use to accelerate learning, as it is capable of expressing manually programmed heuristics like reaching and pulling. Finally, our approach can use the

existing estimated object kinematics to provide a principled reward signal for model-free reinforcement learning.

Recently, residual reinforcement learning approaches have been developed which learn a policy superimposed on hand-designed or model-predictive controllers [37, 14]. Our method is compatible with these approaches, where demonstrations from the motion planner can be used as a base policy on top of which a residual policy can be learned based on kinematic rewards. These methods typically suffer from the same limitations as MPC-based methods mentioned above.

Guided Policy Search (GPS) [23] uses LQR to guide policy search into high-reward regions of the state-space. The models employed are fundamentally local approximations, and thus would benefit greatly from a wealth of suboptimal demonstrations from the outset (as made evident by Chebotar et al. [5]). GPS is one of the state-of-the-art algorithms we expect to be used within our framework as the policy search implementation (Section III-C). A critical distinction between our work and GPS is the notion of planning trajectories in object configuration spaces and reasoning about grasp candidates to achieve a desired manipulation. This is done using information available a priori, and thus is immediately capable of generating high-value policies, whereas GPS is estimating dynamics models given observed data (obtained either from demonstration or random initialization). In the absence of a human demonstrator, our method would provide far more useful data at the outset of learning than running a naively initialized linear-gaussian controller (as evidenced by our comparisons to random initialization). The ideas proposed in our paper are distinct from those put forth in GPS: we present a method for obtaining demonstrations under certain conditions in the absence of a human.

Most similar to our line of work are those that use sample-based motion planners for improved policy learning. Jurgenson and Tamar [15] harness the power of reinforcement learning for neural motion planners by proposing an augmentation of Deep Deterministic Policy Gradient (DDPG) [26] that uses the known robot dynamics to leverage sampling methods like RRT* to reduce variance in the actor update and provide off-policy exploratory behavior for the replay buffer. However, Jurgenson and Tamar [15] are only able to address domains where they can assume good estimates of the dynamics model, such as producing free-space motions to avoid obstacles. Our setting, in contrast, focuses on object manipulation, where dynamics are not readily available, but are critical for learning good policies. Jiang et al. [13] address learning to improve plans produced by a motion planner, but do not bootstrap closed-loop policies. Motion planners aren’t expressive enough to leverage the dynamics in object-manipulation tasks, especially in the presence of unknown dynamics, and traditionally are unable to handle perceptual data like RGB images. Our method, on the other hand, enables motion planning to bootstrap policies that are more expressive than the original planner.

VI. CONCLUSION

We have presented a method that uses kinematic motion planning to bootstrap robot motor policies. By assuming access to a potentially noisy description of the object kinematics, we are able to autonomously generate initial demonstrations that perform as well as human demonstrations, but do not require a human, resulting in a practical method for autonomous motor skill learning.

Our methodology is agnostic to the motion planner, motor policy class, and policy search algorithm, making it a widely applicable paradigm for learning robot motor policies. We demonstrate the power of our methodology by bootstrapping different policy classes with demonstrations from humans and a motion planner, and learn motor policies for three dynamic manipulation tasks: closing a microwave door, opening a drawer, and hitting a ball off a tee. Our framework is the first to enable robots to autonomously bootstrap and improve motor policies with model-free reinforcement learning using only a partially-known kinematic model of the environment.

ACKNOWLEDGEMENT

This research was supported by NSF CAREER Award 1844960 to Konidaris, and by the ONR under the PERISCOPE MURI Contract N00014-17-1-2699. Disclosure: George Konidaris is the Chief Robotist of Realtime Robotics, a robotics company that produces a specialized motion planning processor.

REFERENCES

- [1] Ben Abbatematteo, Stefanie Tellex, and George Konidaris. Learning to generalize kinematic models to novel objects. *Proceedings of the 3rd Conference on Robot Learning*, 2019.
- [2] Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 57(5):469–483, 2009.
- [3] Christopher G Atkeson and Stefan Schaal. Robot learning from demonstration. In *ICML*, volume 97, pages 12–20. Citeseer, 1997.
- [4] Alexandre Campeau-Lecours, Hugo Lamontagne, Simon Latour, Philippe Fauteux, Véronique Maheu, François Boucher, Charles Deguire, and Louis-Joseph Caron L’Ecuyer. Kinova modular robot arms for service robotics applications. In *Rapid Automation: Concepts, Methodologies, Tools, and Applications*, pages 693–719. IGI Global, 2019.
- [5] Yevgen Chebotar, Mrinal Kalakrishnan, Ali Yahya, Adrian Li, Stefan Schaal, and Sergey Levine. Path integral guided policy search. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 3381–3388. IEEE, 2017.
- [6] Ching-An Cheng, Xinyan Yan, Nolan Wagener, and Byron Boots. Fast Policy Learning through Imitation and Reinforcement. *arXiv e-prints*, art. arXiv:1805.10413, May 2018.

- [7] Sachin Chitta, Ioan Sucan, and Steve Cousins. Moveit![ros topics]. *IEEE Robotics & Automation Magazine*, 19(1):18–19, 2012.
- [8] Erwin Coumans et al. Bullet physics library. *Open source: bulletpysics.org*, 15(49):5, 2013.
- [9] Marc Peter Deisenroth, Gerhard Neumann, Jan Peters, et al. A survey on policy search for robotics. *Foundations and Trends® in Robotics*, 2(1–2):1–142, 2013.
- [10] Marcus Gualtieri, Andreas Ten Pas, Kate Saenko, and Robert Platt. High precision grasp pose detection in dense clutter. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 598–605. IEEE, 2016.
- [11] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Advances in neural information processing systems*, pages 4565–4573, 2016.
- [12] Auke J Ijspeert, Jun Nakanishi, and Stefan Schaal. Learning attractor landscapes for learning motor primitives. In *Advances in neural information processing systems*, pages 1547–1554, 2003.
- [13] Yuqian Jiang, Fangkai Yang, Shiqi Zhang, and Peter Stone. Task-motion planning with reinforcement learning for adaptable mobile service robots. In *IROS*, pages 7529–7534, 2019.
- [14] Tobias Johannink, Shikhar Bahl, Ashvin Nair, Jianlan Luo, Avinash Kumar, Matthias Loskyll, Juan Aparicio Ojea, Eugen Solowjow, and Sergey Levine. Residual reinforcement learning for robot control. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 6023–6029. IEEE, 2019.
- [15] Tom Jurgenson and Aviv Tamar. Harnessing reinforcement learning for neural motion planning. *arXiv preprint arXiv:1906.00214*, 2019.
- [16] Gregory Kahn, Tianhao Zhang, Sergey Levine, and Pieter Abbeel. Plato: Policy learning using adaptive trajectory optimization. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3342–3349. IEEE, 2017.
- [17] Sham M Kakade. A natural policy gradient. In *Advances in neural information processing systems*, pages 1531–1538, 2002.
- [18] Bingyi Kang, Zequn Jie, and Jiashi Feng. Policy optimization with demonstrations. In *International Conference on Machine Learning*, pages 2469–2478, 2018.
- [19] Sertac Karaman, Matthew R Walter, Alejandro Perez, Emilio Frazzoli, and Seth Teller. Anytime motion planning using the rrt. In *2011 IEEE International Conference on Robotics and Automation*, pages 1478–1483. IEEE, 2011.
- [20] Oliver Kroemer, Scott Niekum, and George Konidaris. A review of robot learning for manipulation: Challenges, representations, and algorithms. *arXiv preprint arXiv:1907.03146*, 2019.
- [21] Andrey Kurenkov, Ajay Mandlekar, Roberto Martin-

- Martin, Silvio Savarese, and Animesh Garg. Ac-teach: A bayesian actor-critic method for policy learning with an ensemble of suboptimal teachers. *arXiv preprint arXiv:1909.04121*, 2019.
- [22] Steven M LaValle. *Planning algorithms*. Cambridge university press, 2006.
- [23] Sergey Levine and Vladlen Koltun. Guided policy search. In *International Conference on Machine Learning*, pages 1–9, 2013.
- [24] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(39):1–40, 2016. URL <http://jmlr.org/papers/v17/15-522.html>.
- [25] Xiaolong Li, He Wang, Li Yi, Leonidas Guibas, A Lynn Abbott, and Shuran Song. Category-level articulated object pose estimation. *arXiv preprint arXiv:1912.11913*, 2019.
- [26] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [27] Tomas Lozano-Perez. Automatic planning of manipulator transfer movements. *IEEE Transactions on Systems, Man, and Cybernetics*, 11(10):681–698, 1981.
- [28] Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Aparicio Ojea, and Ken Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. *arXiv preprint arXiv:1703.09312*, 2017.
- [29] Ashvin Nair, Dian Chen, Pulkit Agrawal, Phillip Isola, Pieter Abbeel, Jitendra Malik, and Sergey Levine. Combining self-supervised learning and imitation for vision-based rope manipulation. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2146–2153. IEEE, 2017.
- [30] Yunpeng Pan, Ching-An Cheng, Kamil Saigol, Keuntaek Lee, Xinyan Yan, Evangelos Theodorou, and Byron Boots. Agile autonomous driving using end-to-end deep imitation learning. *arXiv preprint arXiv:1709.07174*, 2017.
- [31] Peter Pastor, Heiko Hoffmann, Tamim Asfour, and Stefan Schaal. Learning and generalization of motor skills by learning from demonstration. In *2009 IEEE International Conference on Robotics and Automation*, pages 763–768. IEEE, 2009.
- [32] Jan Peters and Stefan Schaal. Natural actor-critic. *Neurocomputing*, 71(7-9):1180–1190, 2008.
- [33] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- [34] Harish Ravichandar, Athanasios S Polydoros, Sonia Chernova, and Aude Billard. Recent advances in robot learning from demonstration. *Annual Review of Control, Robotics, and Autonomous Systems*, 2019.
- [35] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635, 2011.
- [36] Stefan Schaal and Christopher G Atkeson. Constructive incremental learning from only local information. *Neural computation*, 10(8):2047–2084, 1998.
- [37] Tom Silver, Kelsey Allen, Josh Tenenbaum, and Leslie Kaelbling. Residual policy learning. *arXiv preprint arXiv:1812.06298*, 2018.
- [38] Freek Stulp and Olivier Sigaud. Path integral policy improvement with covariance matrix adaptation. *arXiv preprint arXiv:1206.4621*, 2012.
- [39] Wen Sun, Arun Venkatraman, Geoffrey J Gordon, Byron Boots, and J Andrew Bagnell. Deeply aggravated: Differentiable imitation learning for sequential prediction. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3309–3318. JMLR. org, 2017.
- [40] Tarik Tosun, Eric Mitchell, Ben Eisner, Jinwook Huh, Boram Lee, Daewon Lee, Volkan Isler, H Sebastian Seung, and Daniel Lee. Pixels to plans: Learning non-prehensile manipulation by imitating a planner. *arXiv preprint arXiv:1904.03260*, 2019.
- [41] Mel Vecerik, Todd Hester, Jonathan Scholz, Fumin Wang, Olivier Pietquin, Bilal Piot, Nicolas Heess, Thomas Rothörl, Thomas Lampe, and Martin Riedmiller. Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards. *arXiv preprint arXiv:1707.08817*, 2017.
- [42] Yuke Zhu, Ziyu Wang, Josh Merel, Andrei Rusu, Tom Erez, Serkan Cabi, Saran Tunyasuvunakool, János Kramár, Raia Hadsell, Nando de Freitas, and Nicolas Heess. Reinforcement and imitation learning for diverse visuomotor skills. In *Proceedings of Robotics: Science and Systems*, Pittsburgh, Pennsylvania, June 2018. doi: 10.15607/RSS.2018.XIV.009.