

# Cascaded Compositional Residual Learning for Complex Interactive Behaviors

K. Niranjan Kumar , Irfan Essa , *Fellow, IEEE*, and Sehoon Ha 

**Abstract**—Real-world autonomous missions often require rich interaction with nearby objects, such as doors or switches, along with effective navigation. However, such complex behaviors are difficult to learn because they involve both high-level planning and low-level motor control. We present a novel framework, Cascaded Compositional Residual Learning (CCRL), which learns composite skills by recursively leveraging a library of previously learned control policies. Our framework combines multiple levels of pre-learned skills by using multiplicative skill composition and residual action learning. We also introduce a goal synthesis network and an observation selector to support combination of heterogeneous skills, each with its unique goals and observation space. Finally, we develop *residual regularization* for learning policies that solve a new task, while preserving the style of the motion enforced by the skill library. We show that our framework learns joint-level control policies for a diverse set of motor skills ranging from basic locomotion to complex interactive navigation, including navigating around obstacles, pushing objects, crawling under a table, pushing a door open with its leg, and holding it open while walking through it. The proposed CCRL framework leads to policies with consistent styles and lower joint torques, and successfully transfer to a real Unitree A1 robot without any additional fine-tuning.

**Index Terms**—Legged robots, reinforcement learning, mobile manipulation.

## I. INTRODUCTION

**A**UTONOMOUS missions often involve various levels of motor skills ranging from simple locomotion and manipulation to interactive navigation. For instance, a quadrupedal tasked to fetch an object may need to walk to a nearby door, open it with its end-effector, and navigate to the destination (see Fig. 1). Traditionally, many researchers [1], [2], [3] have approached modeling such high-level behaviors by manually decomposing them into several low-level motor skills, such as locomotion, navigation, manipulation, and a high-level decision layer to modulate these skills. While effective, this model-based approach requires researchers to derive the explicit model of



Fig. 1. Our framework allows robots to learn complex environment interactions with a recursive hierarchy. The image shows a motion frame of the robot crawling under a table, reaching the door, pushing it open, and walking through it to reach a target location.

environmental interactions that are often complicated and cumbersome in such scenarios.

Learning-based approaches, such as deep reinforcement learning (deep RL) [4], [5], [6], hold promise for automatically generating effective motor policies from simple descriptions. However, learning complex behaviors is not straightforward due to many theoretical and practical challenges. For instance, motor creatures with many degrees of freedom (DoF), such as quadrupedal or bipedal robots, require a massive amount of simulation samples even for the simplest motor tasks due to their high-dimensional state and action spaces [7], [8], [9]. In addition, learning a high-level skill requires careful reward engineering, which is extremely time-consuming to tune. Researchers have investigated an explicit two-level hierarchical architecture to alleviate this problem [10], [11], but it only arranges a sequence from existing behaviors. For instance, it cannot combine locomotion and target-reaching into a door-opening skill.

We present Cascaded Compositional Residual Learning (CCRL), a novel framework to learn a family of motor skills on a 12 DoF quadruped robot and interactively navigate around an indoor scene. Our architecture builds new skills by recursively leveraging a pre-learned library consisting of primitive (e.g., walking straight) or other high-level motor skills (e.g., door opening). We accomplish this nested skill composition by learning *multiple* levels of weighting networks (via multiplicative composition [12]) and residual policy networks [13], [14]. Additionally, we introduce *residual regularization* to strike a balance between maximizing a given reward function and adhering to a combination of foundational skills to ensure feasible motion on a real robot.

The main contributions of our work are as follows:

Manuscript received 5 April 2023; accepted 27 May 2023. Date of publication 14 June 2023; date of current version 23 June 2023. This letter was recommended for publication by Associate Editor K. Akbari Hamed and Editor A. Kheddar upon evaluation of the reviewers' comments. This work was supported by the Cisco Research Corporation. (*Corresponding author: K. Niranjan Kumar.*)

K. Niranjan Kumar is with the Georgia Institute of Technology, Atlanta, GA 30332 USA (e-mail: niranjankumar@gatech.edu).

Irfan Essa and Sehoon Ha are with the Georgia Institute of Technology, Atlanta, GA 30332 USA, and also with the Google, Mountain View, CA 94043 USA (e-mail: irfan@computer.org; sehoonha@gatech.edu).

See videos at <https://www.kniranjankumar.com/ccrl/>.

This letter has supplementary downloadable material available at <https://doi.org/10.1109/LRA.2023.3286171>, provided by the authors.

Digital Object Identifier 10.1109/LRA.2023.3286171

- 1) A cascaded compositional residual learning (CCRL) framework to recursively train skills with increasing complexity by reusing skills learned in previous steps.
- 2) A novel combination of multiplicative skill composition and residual learning which serves as the fundamental unit of our multi-level cascaded policy architecture.
- 3) Introduction of a goal synthesis network and an observation selector which allow seamless combinations of heterogeneous skills.
- 4) A constrained residual learning objective that anchors newly acquired skills to achieve realistic and achievable motions.
- 5) Demonstration of interaction skills learned in simulation and successful transfer to a real Unitree A1 robot.

## II. RELATED WORK

### A. Robotic Locomotion

Controlling legged robots has been a long-standing topic of research in the robotics community. Traditionally, roboticists [15], [16], [17], [18], [19], [20] have approached this problem with a combination of different frameworks and algorithms like trajectory optimization, model-predictive and whole-body control to demonstrate robust and agile locomotion. However, such manual controller design techniques often require in-depth prior knowledge about the robot dynamics, which restricts its applicability to new problems. On the other hand, reinforcement learning [5] offers an automated controller design process by optimizing a policy for a reward function that measures performance on a given task. A body of work [7], [8], [9], [21], [22], [23], [24] in RL has demonstrated effective learning of locomotion policies. However, these learned policies often show degraded performance on hardware due to the difference between the simulation and the real world, referred to as the *sim-to-real* gap. Several techniques have been proposed to bridge the sim-to-real gap, such as domain randomization [7], [8], learning actuator dynamics [23], online adaptation [9], [25], real-world learning [26], [27]. Our work also leverages existing domain randomization techniques to deploy learned interactive behaviors to the real world.

### B. Interactive Navigation

Navigation is a fundamental skill for autonomous robot missions. Several navigation problems have been proposed over the years [11], [28], such as: PointNav - navigating to a point in a map; ObjectNav - navigating to a selected object category in the scene; InteractiveNav - Navigating to a point in a scene that requires interaction with objects and furniture in the scene. While researchers have made significant progress in PointNav [29] and ObjectNav [30], [31], [32], InteractiveNav [11], [33], [34], [35] is still a challenging problem due to the difficulty in learning interaction dynamics between a robot and its environment. Li [11] learned a navigation policy on a mobile robot that can open a door to reach its target location. Konidaris [35] took an alternate approach where the robot was given access to a set of hard-coded controllers and had to learn when to use each skill.

However, these approaches focus on just wheeled mobile robots, abstracted to be a simple cylinder with an attached manipulator, significantly limiting the space of possible interactions. Sunwoo [6] demonstrated interaction skills on a quadrupedal robot via manual motion-based control. In contrast, our work focuses on learning autonomous policies to directly control joint motors of a 12-DOF quadruped robot, to perform complex dynamic interactions such as door opening, object pushing and crawling under a table. To the best of our knowledge, our work is the first of its kind, demonstrating end-to-end neural network policies, that can solve InteractiveNav on a high DOF legged robot, with a high success rate.

Another related problem is interactive search, where a robot has to *search* for an object in an environment. Most of the work in this area [36], [37], [38] focuses on searching for an object in a cluttered shelf or table. In a more generalized version of interactive search, a robot moves around a cluttered indoor scene searching for an object. While we do not tackle interactive search, our method complements existing techniques by enabling a robot *interactively* navigate to the object, once its location has been estimated.

### C. Hierarchical Reinforcement Learning

Learning RL policies for complex tasks requires extensive reward engineering and hyperparameter tuning. A common technique to deal with such complexity is to decompose the policy into multiple sub-tasks. In hierarchical RL (HRL) [39], [40], higher level policies control and instruct low level policies like they were primitive actions. While an in-depth discussion about hierarchical RL is outside the scope of this letter, we point interested readers to refer [41], and restrict our discussion to related work in locomotion and navigation. A common approach to navigation is to use low-level locomotion primitives in combination with a high-level primitive selection or goal proposal network to perform navigation [42], [43], [44], [45]. Alternatively, Yang [10] proposed a gating network that combines the parameters of a collection of neural network primitives, to produce composite navigation policy. Similar to these works, we are interested in designing a hierarchy to solve a complex task, but our low-level policies are themselves complex (like opening a door) and not easily breakable into simpler policies. Therefore, a common “flat” hierarchy will not be sufficient to learn more challenging motor skills. We propose a multi-level cascaded hierarchy to build progressively complex policies starting out from a small set of learned skills. We build this hierarchy on the fly, by learning residuals [13], [14] that can perturb the output of an existing library of policies, to generate novel behaviors.

## III. CASCADED COMPOSITIONAL RESIDUAL LEARNING

This section describes our cascaded compositional residual learning (CCRL) framework that learns complex motor behaviors, by recursively obtaining control policies and reusing them to solve increasingly difficult tasks. We will begin the section by explaining relevant background, followed by problem formulation, skill reusing mechanisms and learning algorithms.

### A. Background: Markov Decision Process

Robot learning can be modeled as a Markov Decision Process (MDP) defined by the 5-tuple  $\mathcal{M} : \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, r, \gamma \rangle$ , where  $\mathcal{S}$ ,  $\mathcal{A}$  and  $r$  are the state space, the action space, and the reward function, respectively.  $\mathcal{T}$  is a transition function that determines the next state given the current state and action.  $\gamma \in (0, 1)$  is the discount factor. Our goal then, is to find a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  that maximizes a cumulative return  $J(\pi) = \mathbb{E}_{\tau \sim \rho_\pi} [\sum_{t=0}^T r(s_t, a_t)]$ , where  $\rho_\pi$  is the distribution over trajectories ( $\tau$ ) imposed by policy  $\pi$  and  $\mathbb{E}$  is the expectation.

However, solving the given MDPs for high-level motor tasks, such as interactive navigation or full-body manipulation, is not straightforward. First, the given behavior involves multiple sub-tasks, such as walking or door opening, which require the manual design of very specific reward functions. In addition, the long-horizon nature of the problem makes it extremely sensitive to the choice of hyper-parameters. Further, we want to maintain smooth and stable motion styles to make them consistently feasible on actual hardware. Therefore, it is almost impossible to obtain such good motion controllers for these challenging motor tasks via simple reward engineering and hyperparameter tuning.

### B. Background: Residual Learning

Residual learning [13], [14] offers a promising alternative to the traditional training-from-scratch RL framework by enabling the transfer of a skill learned on one task to another related but more difficult task. Consider a policy  $\pi_0$  to be a solution for the base task that takes as input, state  $s$  and outputs action  $a$ . We can freeze the parameters of  $\pi_0$  and learn a residual policy  $\pi_k$  to perturb the actions of  $\pi_0$ , to modify learned behaviors and accomplish a different goal. The final policy is obtained by adding the outputs:

$$\pi(s) = \pi_0(s) + \pi_k(s) \quad (1)$$

### C. Problem Definition: Multi-Skill Learning

In our formulation, we consider a set of  $n$  motor tasks described as MDPs  $\Omega = \{\mathcal{M}_0, \mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_n\}$ , where each  $\mathcal{M}_k : \langle \mathcal{S}_k, \mathcal{A}, \mathcal{T}, r_k, \gamma \rangle \quad \forall \quad 0 < k < n$ , models the MDP for a learnable skill.  $\mathcal{S}_k, \mathcal{A}, \mathcal{T}, r_k$  and  $\gamma$  represent the state space, action space, transition function, reward function and discount factor respectively. Note that we allow unique state spaces for the problems, which indicates that some high-level behaviors require additional inputs for describing the task (e.g., the target location) or environmental states (e.g., the door hinge angle). On the other hand, we assume a single robot, which leads to unified actions and transition functions.

We then want to learn a collection of control policies  $\mathcal{C} = \{\pi_1, \pi_2, \pi_3, \dots, \pi_n\}$ . A naïve approach is to independently solve each corresponding MDP. But in practice, it is often not feasible due to the challenges outlined in Section III-A. Our key insight is to “recursively” learn policies for challenging MDPs on top of the related “prerequisite” or “parent” skills instead of learning from scratch. To this end, we assume a Skill Decomposition Graph  $\mathcal{G}$  by defining a set of prerequisite skills, where its vertices are the corresponding MDPs, and its

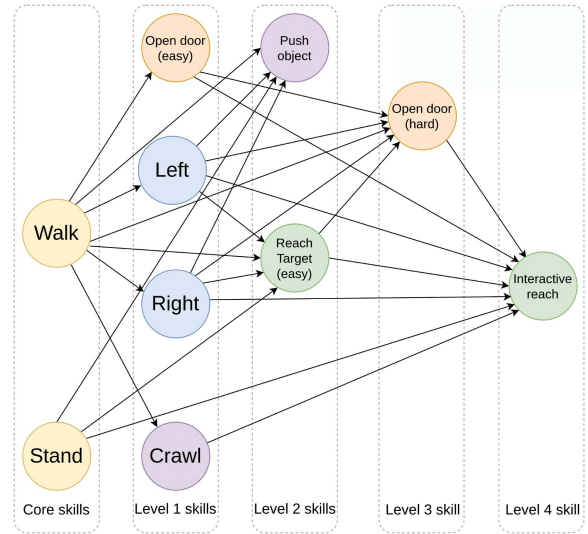


Fig. 2. Directed graph representing the relationship between skills. Each skill is built as a compositional policy over policies from earlier levels and a learned residual.

edges represent the dependencies. For instance, Fig. 2, shows the dependencies of *Reach Target (Easy)* on *Walk*, *Stand*, *Turn Left* and *Turn Right*. We assume the skill decomposition graph  $\mathcal{G}$  to be intuitive to define (e.g., navigation requires straight walking and turning) but could also be automatically estimated by searching the space of all possible relations and retaining those that yield the best performance. Further, our framework is robust to the redundancy in prerequisite skills in the sense that it learns with subset of them or ignore irrelevant prerequisite skills (e.g., if the task does not require the robot to crawl, the crawling skill will be ignored). Please refer to Section IV-D.

### D. Cascaded Compositional Residual Learning

A traditional flat architecture of HRL would train each of these policies from scratch, to solve the corresponding MDP. But due to challenges outlined in the previous section, it is not practically feasible. The key insight of our approach is that, instead of learning new skills from scratch, we build them on top of pre-learned skills, i.e.  $\pi_k$  could be a composite policy over the prerequisite policies  $\mathcal{C}_k$ , where  $\mathcal{C}_k \subset \mathcal{C}$ . Mathematically,

$$\pi_k = \mathcal{F}_k(\mathcal{C}_k \cup \{\pi_k^r\}) \quad (2)$$

where  $\mathcal{F}_k$  is a merging function that combines the outputs of all  $\pi_j \quad \forall \quad \pi_j \in \mathcal{C}_k$  and  $\pi_k^r$  is a residual neural network policy. To this end, our framework (Fig. 3) has three trainable networks, (1) a residual action network, (2) a weight network, and (3) a synthetic goal network, while freezing the prerequisite policies in  $\mathcal{C}_k$ .

*Residual Network:* The residual policy  $\pi_k^r$  generates a *residual action* that perturbs the actions generated by the learned policies to generate novel behavior that is absent in the set of pre-learned skills. We simply treat this residual network as an additional policy: the only difference being that the parameters of  $\pi_k^r$  are learnable during policy training.

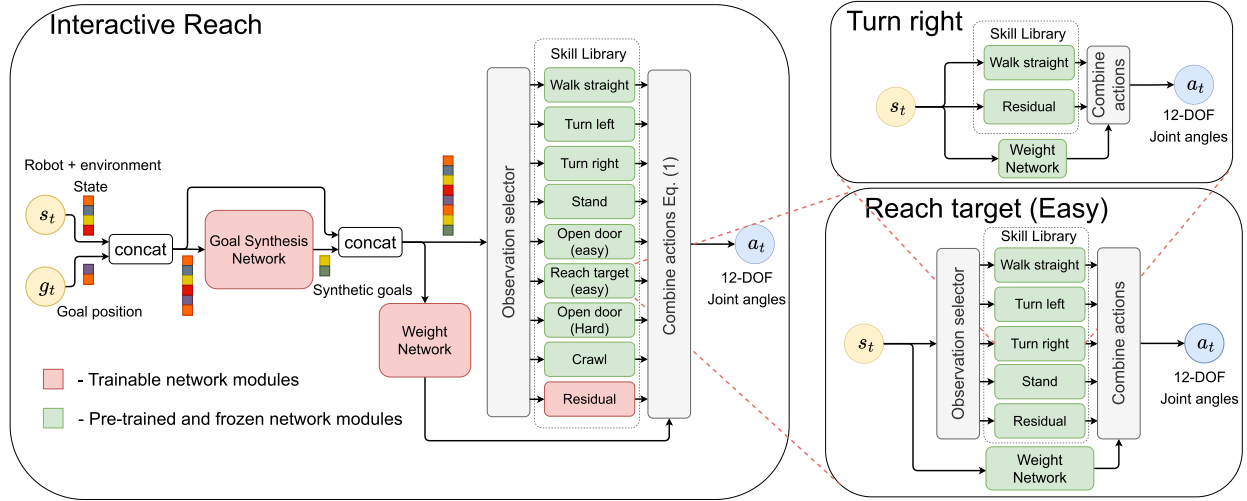


Fig. 3. Overview of our Cascaded Compositional Residual Learning, which consists of three learnable components: a synthetic goal network, a weight network, and a residual action network. The example illustrates the architecture for *Interactive Reach*, which uses *Reach target (Easy)* on top of other primitive skills.

*Weight Network*: While prior work in residual learning combines two policies by just adding them together, this approach has been observed to be sub-optimal by Peng et al. [12] when combining a large set of skills. They proposed Multiplicative Compositional Policy (MCP) that combines a set of stochastic policies while ensuring that multiple policies simultaneously work together and provide more flexibility in the final composite policy. MCP combines multiple skills by using the following equation:

$$\pi(a | s, g) = \frac{1}{Z(s, g)} \prod_{i=0}^k \pi_i(a | s, g)^{w_i(s, g)}, \quad (3)$$

where  $w_i(s, g)$  are learned weighting functions and  $Z(s, g)$  is a normalizing factor to ensure that the weights are between 0 and 1. We use MCP as the merging function  $\mathcal{F}_k$ .

*Goal Synthesis Network*: It is not intuitive to combine cascaded motor skills because some skills take heterogeneous goals as additional inputs. For example, the *Reach Target (Easy)* skill needs a goal location around the robot as input. But to use this skill as part of *Door open (Hard)*, we need to generate a collection of intermediate targets as input the policy. These targets are generated by the Synthetic Goals Network. We normalize the network outputs, scale them back to feasible goal ranges and concatenate them to the observations.

*Observation Selector*: Since each skill has its own observation space, we introduce an *Observation selector* which takes the input observations from the environment and selects the subset of observations relevant to each skill. The architecture is illustrated in Fig. 3. Our work provides a framework that builds novel complex skills, by leveraging a library of pre-learned composite skills, while ensuring control over the style of the composite. We use our framework to learn a wide range of interactive skills on a quadruped robot while implicitly guiding the *style* of the policy by controlling the weight given to the residual actions. With this framework we train a handful of basic policies from

scratch, using RL, and use these policies to progressively build increasingly complex composite skills.

### E. Auxiliary Loss for Constrained Residual Actions

In the previous section, we discussed learning novel skills by effectively leveraging a library of existing skills and a trainable residual. However, an unconstrained residual could dominate the output of the policy and give rise to unstable behaviors, which eventually lead to sub-optimal behaviors and unsuccessful sim-to-real transfer. Therefore, we introduce constraints into our policy optimization framework to control the extent to which residuals influence the policy. In our work, we use PPO (Proximal Policy Optimization) [46] as the RL algorithm of choice. Consider the standard PPO objective below:

$$L_t^{\text{PPO}}(\theta) = \hat{\mathbb{E}}_t [L_t^{\text{CLIP}}(\theta) - c_1 L_t^{\text{VF}}(\theta) + c_2 S[\pi_\theta](s_t)] \quad (4)$$

where  $L_t^{\text{CLIP}}(\theta)$  is the clipped policy gradient objective,  $L_t^{\text{VF}}(\theta)$  is the value function error and  $S[\pi_\theta](s_t)$  is the entropy term.  $c_1, c_2$  are weights to control the influence of each term over the total loss. We set coefficients  $c_1, c_2$  to be 1.0 and 0.01 respectively. To constrain the effect of the residual, we introduce regularization terms to penalize the magnitude and weight ( $w_{\text{res}}$ ) given to the residual.

$$L_t^{\text{rm}} = | \mathbb{E}_t [\pi_k^r(a_t | s_t, g_t)] |_1 \quad (5)$$

$$L_t^{\text{rw}} = \mathbb{E}_t [| w_{\text{res}}(s_t, g_t) |_1] \quad (6)$$

$L_t^{\text{rm}}$  and  $L_t^{\text{rw}}$  are penalties on the magnitude of the residual actions and the weights assigned to them respectively.

We modify PPO objective (4) to include regularization terms defined above:

$$L_t(\theta) = L_t^{\text{PPO}}(\theta) + c_3 L_t^{\text{rm}} + c_4 L_t^{\text{rw}}. \quad (7)$$

$c_3$ , and  $c_4$  are task dependent coefficients that are tuned to maximize the policy performance (see Table I). Once we train a skill, we save the residual, weight, and goal synthesis networks

TABLE I  
NETWORK ARCHITECTURES

Skill	Residual Architecture	Weight Net. Architecture	$c_3$	$c_4$
Walk straight	[512, 256, 128]	-	-	-
Stand	[256, 128]	-	-	-
Turn left	[256, 128]	[128]	0.01	0.01
Turn right	[256, 128]	[128]	0.01	0.01
Reach target(Easy)	[256, 128]	[512, 256, 128]	0.05	0
Open door(Easy)	[512, 256]	[512, 256, 128]	0.01	0.01
Open door(Hard)	[512, 256]	[256, 128]	0.9	0
Push object	[256, 128]	[512, 256, 128]	0.5	0.5
Crawl	[256, 128]	[512, 256, 2]	0.03	0
Interactive Reach	[256, 128]	[512, 256, 2]	5.0	5.0

and reuse them to train other skills. The critic is trained from scratch for every new skill.

#### IV. EXPERIMENTS

In this section, we design experiments to test the effectiveness of cascaded compositional residual learning (CCRL), by learning skills involving different levels of difficulty and comparing the performance with traditional RL baselines. Next, we present simulation analysis to investigate the sensitivity of the performance with respect to the skill choices and the importance of the auxiliary loss. Finally, we examine the learned skill policies on a real Unitree A1 robot and demonstrate robust sim-to-real transfer.

##### A. Problem Formulation

In our work, we train a collection of policies, that help an agent interactively navigate an indoor scene. Each skill has its own underlying MDP, but we provide a rough definition of the ingredients of the MDPs as follows.

*States:* A 60-dimensional state space consists of the base linear velocity (3), base angular velocity (3), base yaw angle (1), base position (2), gravity vector (3) projected on the base, joint angles (12), joint velocities (12), previous action (12), locations  $(x, y)$  of objects and the target relative to the robot (8), angle of the door joint(1), room index of the robot and target (2) and crawling height (1).

*Actions:* The action space is the desired joint angles of the robot (12 for A1), which are fed into Proportional Derivative (PD) controllers.

*Rewards:* We use a collection of terms to account for different aspects of the robot motion. Let us denote the velocity as  $\mathbf{v}$ , the angular velocity as  $\boldsymbol{\omega}$ , the joint angles as  $\mathbf{q}$ , joint velocities as  $\dot{\mathbf{q}}$ , joint torques as  $\boldsymbol{\tau}$ , number of robot parts (excluding feet) in contact with the environment as  $n_{\text{contact}}$ , action taken at a given step as  $a_t$ , the angle made by the hinge of the door joint as  $q_{\text{door}}$ , the position of the target relative to the robot as  $\mathbf{x}_{\text{target}}$ , the position of the target relative to the object being pushed as  $\mathbf{x}_{t2o}$  and  $\Delta t$  as the simulation time-step. The reward at time  $t$  is defined as the weighted sum of the following quantities:

- R1) Linear velocity tracking:  $\exp(-(\mathbf{v}_{\text{target}} - \mathbf{v})^2/\sigma_1)$
- R2) Angular velocity tracking:  $\exp(-(\boldsymbol{\omega}_{z_{\text{target}}} - \boldsymbol{\omega}_z)^2/\sigma_2)$
- R3) Pitch and roll penalty:  $\omega_x^2 + \omega_y^2$
- R4) Joint acceleration penalty:  $\sum (\frac{\dot{\mathbf{q}}_t - \dot{\mathbf{q}}_{t-1}}{\Delta t})^2$
- R5) Collision penalty:  $n_{\text{contact}}$

R6) Action change penalty:  $\sum (\mathbf{a}_t - \mathbf{a}_{t-1})^2$

R7) Torque penalty:  $\sum \boldsymbol{\tau}^2$

R8) Door angle:  $q_{\text{door}}$

R9) Distance to target:  $\exp(-\|\mathbf{x}_{\text{target}}\|/\sigma_3)$

R10) Object-target distance:  $\exp(-\|\mathbf{x}_{t2o}\|/\sigma_3)$

$\sigma_1, \sigma_2, \sigma_3$  are scaling factors that we set to 0.25, 0.25 and 2.0 respectively.

*Tasks:* We define the following tasks, each aimed at developing a skill that is useful while navigating or exploring an indoor cluttered scene (Fig. 4). Their relationship is defined by the skill graph  $\mathcal{G}$  (Fig. 2).

- 1) **Walk straight:** Walk forward in a straight line.
- 2) **Stand:** Stand in place when subjected to external disturbances.
- 3) **Turn left:** Turn to the left. The episode is successful if the robot circle backs to its initial yaw.
- 4) **Turn right:** Turn to the right. The episode is successful if the robot circle backs to its initial yaw.
- 5) **Reach target (Easy):** Reach a target placed anywhere within a circle of radius 3 m from the center of the robot. The episode is successful if the robot's center is within 0.5 m from the goal (R9).
- 6) **Open door (Easy):** Push door open with foot and hold it open while walking through it. The episode is successful if the robot crosses the door (R8).
- 7) **Open door (Hard):** Reach the door from an arbitrary starting point in a room, push it open with a foot and hold it open while walking through it. The episode is successful if the robot crosses the door (R8).
- 8) **Push object:** Push a cylindrical puck to a target location. The episode is successful if the puck is within 0.1 m from the desired target location (R10).
- 9) **Crawl:** Crawl under a slab that gradually decreases in height, while avoiding collisions. The episode is successful if the robot crawls out from under the slab.
- 10) **Interactive Reach:** Reach a target object placed anywhere in a two-room house with the fixed layout of furniture and a door, while avoiding collisions. Robot is successful if its body center within 0.5 m from the center of the target (R9).

All the tasks are trained with rewards R1-R7 in addition to the ones mentioned above.

##### B. Simulation Setup

We use Isaac Gym [47] to simulate our interactive environment and train our policies. We run 4096 environments in parallel, on a single NVIDIA Titan X GPU. During training, we randomize surface friction between the robot and the ground by randomly sampling from the range-[0.5, 1.25]. We implement all our policies as fully connected neural network layers (see Table I) with exponential linear units (*Elu*). We train policies with our updated PPO objective (7).

##### C. Simulation Results

Our framework enables us to train complex long horizon interactive behaviors grounded in previously learned parent skills. Our policy was able to learn all the tasks described in

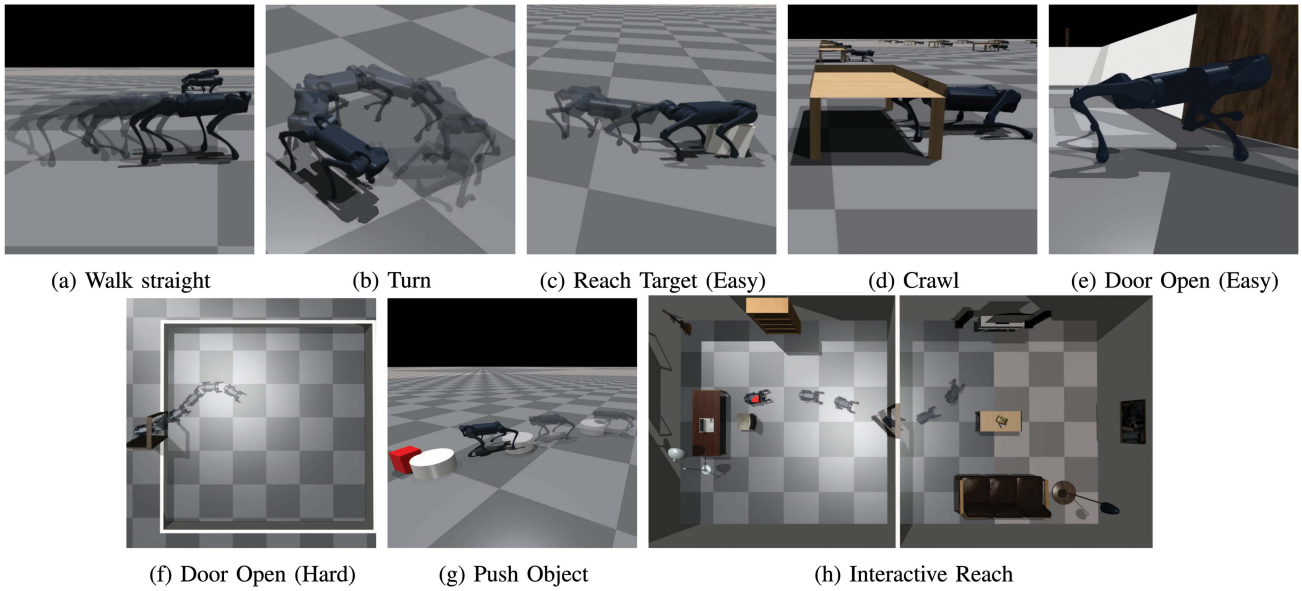


Fig. 4. Overview of the tasks we train our robot to accomplish.

Section IV-A and Fig. 4. The success rates (over 100 trials) of the learned policies are near 90%, except for the hardest task of *Interactive Reach* that shows a 72% success rate.

*Performance Comparison:* To highlight the importance of our cascaded residual framework, we compare our method against four baselines:

- 1) **Vanilla:** The policy is trained from scratch to solve the given problem, using the same architecture as the *Walk straight* skill.
- 2) **Curriculum:** A single policy is trained using a manually curriculum that gradually increases the complexity of the environment, using the same architecture as the *Walk straight* skill.
- 3) **Big Policy:** To show that our improved performance is not a byproduct of architecture size, we train a policy that uses the same neural network architecture as our *Interactive Reach* policy including sub-networks, but is trained from scratch.
- 4) **MCP+Residual:** To emphasize the importance of residual penalties, we compare against a naive combination of MCP and Residual learning by removing auxiliary regularization losses.

We train all policies until convergence and employ early stopping when required. To train all the skills in our framework we used a total of 3 billion environment steps. For most of the baselines, we stopped seeing any improvement in the reward or success rate beyond 0.5 billion environment steps for tasks in levels 1–3. For our *Interactive reach* task, we trained all our policies for 1.5 billion steps after which there is no noticeable improvement in performance. We compare the success rates on multiple tasks in Table II. Our approach, CCRL, typically outperforms all the other three baselines, Vanilla, Curriculum, and Big Policy, on challenging tasks. For the *Door Open (Hard)*, only CCRL and Curriculum achieve good success rates (98%) while Vanilla and Big policy fail. In our experience, it is almost

TABLE II  
COMPARISON OF SUCCESS RATES IN SIMULATION

	Success Rate			
	Reach Target (Easy)	Open Door (Hard)	Push Object	Interactive Reach
Vanilla	0.96	0	0.04	0.46
Curriculum	0.92	<b>0.98</b>	0.52	0.48
Big Policy	0.92	0	0.01	0.22
MCP+Residual	<b>0.98</b>	0.12	0	0.48
<b>CCRL</b>	<b>0.98</b>	<b>0.98</b>	<b>0.89</b>	<b>0.72</b>

impossible to learn an effective door-opening skill without a proper curriculum, such as learning to walk and open a door. However, for the most complicated task of *Interactive Reach*, CCRL shows the best success rate of 72% while all three baselines only show less than 48% success rates. This is because *Interactive Reach* requires careful coordination of all multiple different skills, navigation, crawling, obstacle avoidance, and door manipulation, which shows the importance of cascaded skill learning. For rest of the tasks our method reaches a success rate of 1.0 achieving the goal in every episode.

*Control over Quality:* The cascaded architecture also provides a natural way for controlling the style of the motion. For instance, in the case of *Reach Target (easy)*, the robot might jump toward the target or wiggle its legs to move gradually toward the target. While these behaviors might be successful in simulation, they do not transfer well to the real world as they are over-fitted to the simulation dynamics and parameters. To evaluate the quality of trajectories generated by the policy, we use two metrics: Similarity Index (SI) [48] and average torque. The SI is defined as  $SI(X_L, X_R) = 2|X_L - X_R| / (X_L + X_R)$ , where  $X_L$  and  $X_R$  are the average torques used by the left and right motors.

We evaluate these metrics on *Reach Target (Easy)* in Table III for a fixed target location. CCRL uses a more symmetric gait due

TABLE III  
COMPARISON OF TRAJECTORY QUALITY

	Symmetry Index	Average Torque ( $N \cdot m$ )
Vanilla	0.30	6.72
Curriculum	0.26	7.21
Big Policy	0.29	3.33
MCP+Residual	0.28	6.05
<b>CCRL</b>	<b>0.18</b>	<b>2.84</b>

TABLE IV  
SUCCESS RATES ON HARDWARE EXPERIMENTS

	Turn	Crawl	Reach Target	Door Open	Push Object
Curriculum	0.2	0	0	0	0
<b>CCRL</b>	<b>1</b>	<b>1</b>	<b>0.9</b>	<b>0.9</b>	<b>0.8</b>

to its grounding on the core skills through the residual penalties and the lowest average torque when compared to baselines.

#### D. Simulation Analysis

This section aims to analyze the sensitivity of our learning framework, CCRL, by comparing the performance changes against the original formulation.

*Sensitivity to bad parent skills:* We first investigate the robustness of the proposed CCRL to check if it can learn to ignore a bad parent policy with low rewards. To this end, we learn a new policy on the *Reach Target (Easy)* task, but add one more expert policy that is randomly initialized and untrained. We notice that this policy performs similarly to our best policy (CCRL) in Table II, reaching a success rate of 98%. This showcases the robustness of CCRL, in ignoring bad parent policies.

*Sensitivity to irrelevant parent skills:* We also check whether CCRL is sensitive to the design of the skill graph. Particularly, we investigate if CCRL can be robust to additional unnecessary skill dependencies, i.e., additional edges. We retrain a new policy to walk in a straight line with a redundant set of skills:  $\{Walk\ straight, Turn\ left, Turn\ right\}$ .

Our framework ignores the redundancy and walks straight achieving an average velocity tracking error of about 16%, same as the parent walking skill.

*Importance of auxiliary loss:* We further examine the importance of residual penalties ( $L_t^{rm}, L_t^{rw}$ ) by retraining the *Reach Target (Easy)* policy without residual penalties ( $c_3, c_4 = 0$ ). With this setup, the policy ends up overusing residual actions to learn behaviors that diverge too much from the base walking skill. While it still reaches the target in many cases, the quality of trajectories is poor and unsuitable sim-to-real transfer. Please refer to the supplemental video for comparison.

*Importance of residual actions:* The importance of residual actions varies per task. For instance, *Push Object* or *Open Door*, show a zero percent success rate if we ablate the residual actions. On the other hand, *Reach Target* does not rely on residual actions by achieving 96% success that is only 2% lower than the original.

#### E. Hardware Experiments

We transfer our policies trained in simulation to a real Unitree A1 quadruped robot. We use motion capture to identify and track the position of the robot and objects in the real world. Our policies do not require any additional finetuning on a real robot, demonstrating the robustness of the skills generated using CCRL.

In Table IV, we show the success rate of our method compared to our strongest baseline, Curriculum learning. The other

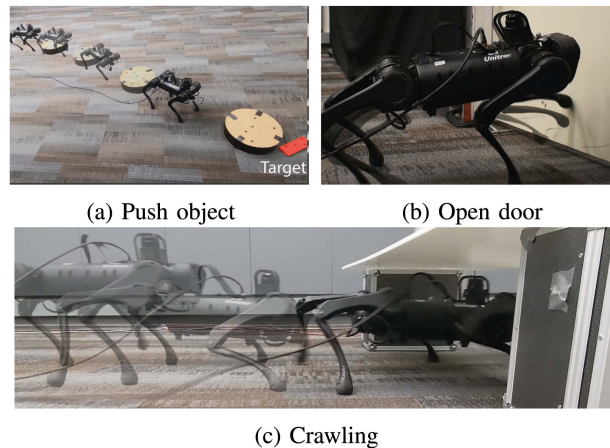


Fig. 5. We transfer the skills, (a) *Push object*, (b) *Open door* and (c) *Crawling* to the real world without any finetuning.

baselines were not successful on hardware. For each skill, we ran every policy 10 times on the real robot. Our method demonstrates robust and reliable sim2real transfer, achieving a high success rate (0.8) on even challenging tasks like *Push object*. Curriculum learning cannot be successfully transferred to even simplest tasks such as *Turn* when subject to the same level of domain randomization, despite reasonable simulation performance. This is a consequence of poor trajectory quality (see Table III) resulting in abrupt and unstable movement on real robot hardware.

In Fig. 1 we show an interactive navigation task, where the robot crawls under a desk, reaches a door, pushes it open with its leg, and walks through it to reach a target. We show example motions of additional skills such as *Push Object*, *Open Door (Easy)* and *Crawling* in Fig. 5.

## V. CONCLUSION

In this letter, we presented a novel Cascaded Compositional Residual Learning (CCRL) framework that recursively learns policies by grounding them to a set of prerequisite skills learned in previous iterations. Using CCRL, we built interactive motor controllers on a high-dimensional quadrupedal robot while ensuring that the learned policies follow a style grounded in the parent skill library. We compared the policies learned to multiple baselines and showed the effectiveness of the proposed framework. We also demonstrated sim-to-real transfer of the learned motor skills.

For future work, we aim to automatically discover skill relationships instead of manually designing the skill graph. Additionally, our current policies adapt the behaviors based on the motion capture and memorization over the fixed layout. We hope to extend our approach to vision-based obstacle avoidance

tasks by learning image-conditioned residual perturbations to pre-trained navigation skills. Additionally, we want to explore interaction policies that learn about objects [49] and adapt to the variability seen in the world.

## REFERENCES

- [1] L. Sentis and O. Khatib, "A whole-body control framework for humanoids operating in human environments," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2006, pp. 2641–2648.
- [2] H. Ferrolho, V. Ivan, W. Merkt, I. Havoutis, and S. Vijayakumar, "Roloma: Robust loco-manipulation for quadruped robots with arms," 2022, *arXiv:2203.01446*.
- [3] S. Zimmermann, R. Poranne, and S. Coros, "Go fetch!-dynamic grasps using boston dynamics spot with external robotic arm," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2021, pp. 4488–4494.
- [4] K. Arulkumar, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep reinforcement learning: A brief survey," *IEEE Signal Process. Mag.*, vol. 34, no. 6, pp. 26–38, Nov. 2017.
- [5] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 2018.
- [6] S. Kim, M. Sorokin, J. Lee, and S. Ha, "Human motion control of quadrupedal robots using deep reinforcement learning," in *Proc. Robot.: Sci. Syst.*, 2022.
- [7] J. Tan et al., "Sim-to-real: Learning agile locomotion for quadruped robots," in *Proc. Robot.: Sci. Syst.*, 2018.
- [8] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning robust perceptive locomotion for quadrupedal robots in the wild," *Sci. Robot.*, vol. 7, no. 62, 2022, Art. no. eabk2822.
- [9] A. Kumar, Z. Fu, D. Pathak, and J. Malik, "RMA: Rapid motor adaptation for legged robots," in *Proc. Robot.: Sci. Syst.*, 2021.
- [10] C. Yang, K. Yuan, Q. Zhu, W. Yu, and Z. Li, "Multi-expert learning of adaptive legged locomotion," *Sci. Robot.*, vol. 5, no. 49, 2020, Art. no. eabb2174.
- [11] C. Li, F. Xia, R. Martin-Martin, and S. Savarese, "HRL4IN: Hierarchical reinforcement learning for interactive navigation with mobile manipulators," in *Proc. Conf. Robot. Learn.*, 2020, pp. 603–616.
- [12] X. B. Peng, M. Chang, G. Zhang, P. Abbeel, and S. Levine, "MCP: Learning composable hierarchical control with multiplicative compositional policies," *NeurIPS*, vol. 32, pp. 3681–3692, 2019.
- [13] T. Silver, K. Allen, J. Tenenbaum, and L. Kaelbling, "Residual policy learning," 2018, *arXiv:1812.06298*.
- [14] T. Johannink et al., "Residual reinforcement learning for robot control," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2019, pp. 6023–6029.
- [15] T. Apgar, P. Clary, K. Green, A. Fern, and J. W. Hurst, "Fast online trajectory optimization for the bipedal robot cassie," in *Proc. Robot.: Sci. Syst.*, vol. 101, Jun. 2018, p. 14.
- [16] M. Hutter et al., "ANYmal-a highly mobile and dynamic quadrupedal robot," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2016, pp. 38–44.
- [17] G. Bledt, M. J. Powell, B. Katz, J. Di Carlo, P. M. Wensing, and S. Kim, "MIT cheetah 3: Design and control of a robust, dynamic quadruped robot," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2018, pp. 2245–2252.
- [18] J. Di Carlo, P. M. Wensing, B. Katz, G. Bledt, and S. Kim, "Dynamic locomotion in the MIT cheetah 3 through convex model-predictive control," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2018, pp. 1–9.
- [19] M. H. Raibert, "Trotting, pacing and bounding by a quadruped robot," *J. Biomech.*, vol. 23, pp. 79–98, 1990.
- [20] D. Kim, J. Di Carlo, B. Katz, G. Bledt, and S. Kim, "Highly dynamic quadruped locomotion via whole-body impulse control and model predictive control," 2019, *arXiv:1909.06586*.
- [21] Z. Fu, A. Kumar, J. Malik, and D. Pathak, "Minimizing energy consumption leads to the emergence of gaits in legged robots," in *Proc. Conf. Robot Learn.*, 2021, pp. 928–937.
- [22] Z. Fu, A. Kumar, A. Agarwal, H. Qi, J. Malik, and D. Pathak, "Coupling vision and proprioception for navigation of legged robots," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops*, 2022, pp. 17273–17283.
- [23] J. Hwangbo et al., "Learning agile and dynamic motor skills for legged robots," *Sci. Robot.*, vol. 4, no. 26, 2019, Art. no. eaau5872.
- [24] N. Rudin, D. Hoeller, P. Reist, and M. Hutter, "Learning to walk in minutes using massively parallel deep reinforcement learning," in *Proc. Conf. Robot Learn.*, 2022, pp. 91–100.
- [25] W. Yu, J. Tan, Y. Bai, E. Coumans, and S. Ha, "Learning fast adaptation with meta strategy optimization," *IEEE Robot. Autom. Lett.*, vol. 5, no. 2, pp. 2950–2957, Apr. 2020.
- [26] S. Ha, P. Xu, Z. Tan, S. Levine, and J. Tan, "Learning to walk in the real world with minimal human effort," in *Proc. Conf. Robot Learn.*, 2020, pp. 1110–1120.
- [27] L. Smith, J. C. Kew, X. B. Peng, S. Ha, J. Tan, and S. Levine, "Legged robots that keep on learning: Fine-tuning locomotion policies in the real world," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2022, pp. 1593–1599.
- [28] P. Anderson et al., "On evaluation of embodied navigation agents," 2018, *arXiv:1807.06757*.
- [29] E. Wijmans et al., "DD-PPO: Learning near-perfect pointgoal navigators from 2.5 billion frames," in *Proc. Int. Conf. Learn. Representations*, 2020.
- [30] M. Sorokin, W. Yu, S. Ha, and C. K. Liu, "Learning human search behavior from egocentric visual inputs," in *Proc. Comput. Graph. Forum*, 2021, pp. 389–398.
- [31] D. S. Chaplot, D. P. Gandhi, A. Gupta, and R. R. Salakhutdinov, "Object goal navigation using goal-oriented semantic exploration," *NeurIPS*, vol. 33, pp. 4247–4258, 2020.
- [32] D. Batra et al., "Objectnav revisited: An evaluation of embodied agents navigating to objects," 2020, *arXiv:2006.13171*.
- [33] F. Xia et al., "Interactive gibbon benchmark: A benchmark for interactive navigation in cluttered environments," *IEEE Robot. Autom. Lett.*, vol. 5, no. 2, pp. 713–720, Apr. 2020.
- [34] K.-H. Zeng, L. Weihs, A. Farhadi, and R. Mottaghi, "Pushing it out of the way: Interactive visual navigation," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2021, Art. no. 9868.
- [35] G. Konidaris, S. Kuindersma, R. Grupen, and A. Barto, "Autonomous skill acquisition on a mobile manipulator," in *Proc. AAAI*, 2011, vol. 25, no. 1, pp. 1468–1473.
- [36] K. N. Kumar, I. Essa, and S. Ha, "Graph-based cluttered scene generation and interactive exploration using deep reinforcement learning," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2022, pp. 7521–7527.
- [37] M. Danielczuk, A. Angelova, V. Vanhoucke, and K. Goldberg, "X-ray: Mechanical search for an occluded object by minimizing support of learned occupancy distributions," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2020, pp. 9577–9584.
- [38] H. Huang et al., "Mechanical search on shelves using lateral access x-ray," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2021, pp. 2045–2052.
- [39] P.-L. Bacon, J. Harb, and D. Precup, "The option-critic architecture," in *Proc. AAAI Conf. Artif. Intell.*, 2017, pp. 1726–1734.
- [40] R. S. Sutton, D. Precup, and S. Singh, "Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning," *Artif. Intell.*, vol. 112, no. 1-2, pp. 181–211, 1999.
- [41] S. Pateria, B. Subagdja, A.-h. Tan, and C. Quek, "Hierarchical reinforcement learning: A comprehensive survey," *ACM Comput. Surv.*, vol. 54, no. 5, pp. 1–35, 2021.
- [42] T. Li, N. Lambert, R. Calandra, F. Meier, and A. Rai, "Learning generalizable locomotion skills with hierarchical reinforcement learning," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2020, pp. 413–419.
- [43] D. Jain, A. Iscen, and K. Caluwaerts, "Hierarchical reinforcement learning for quadruped locomotion," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2019, pp. 7551–7557.
- [44] N. Heess, G. Wayne, Y. Tassa, T. Lillicrap, M. Riedmiller, and D. Silver, "Learning and transfer of modulated locomotor controllers," 2016, *arXiv:1610.05182*.
- [45] X. B. Peng, G. Berseth, K. Yin, and M. Van De Panne, "Deeploco: Dynamic locomotion skills using hierarchical deep reinforcement learning," *Trans. Graph.*, vol. 36, no. 4, pp. 1–13, 2017.
- [46] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," 2017, *arXiv:1707.06347*.
- [47] V. Makovychuk et al., "Isaac gym: High performance GPU based physics simulation for robot learning," in *Proc. NeurIPS Datasets Benchmarks Track*, 2021.
- [48] W. Yu, G. Turk, and C. K. Liu, "Learning symmetric and low-energy locomotion," *ACM Trans. Graph.*, vol. 37, pp. 1–12, 2018.
- [49] K. N. Kumar, I. Essa, S. Ha, and C. K. Liu, "Estimating mass distribution of articulated objects using non-prehensile manipulation," 2019, *arXiv:1907.03964*.