

# Efficient Object Manipulation Planning with Monte Carlo Tree Search

Huaijiang Zhu<sup>1</sup>, Avadesh Meduri<sup>1</sup>, Ludovic Righetti<sup>1</sup>

**Abstract**—This paper presents an efficient approach to object manipulation planning using Monte Carlo Tree Search (MCTS) to find contact sequences and an efficient ADMM-based trajectory optimization algorithm to evaluate the dynamic feasibility of candidate contact sequences. To accelerate MCTS, we propose a methodology to learn a goal-conditioned policy-value network and a feasibility classifier to direct the search towards promising nodes. Further, manipulation-specific heuristics enable to drastically reduce the search space. Systematic object manipulation experiments in a physics simulator and on real hardware demonstrate the efficiency of our approach. In particular, our approach scales favorably for long manipulation sequences thanks to the learned policy-value network, significantly improving planning success rate. All source code including the baseline can be found at <https://hzhu.io/contact-mcts>.

## I. INTRODUCTION

The ability to plan sequences of contacts and movements to manipulate objects is central to endow robots with sufficient autonomy to perform complex tasks. This remains, however, particularly challenging. Indeed, finding dynamically feasible sequences of contacts between the manipulator and an object typically leads to intractable combinatorial and nonlinear problems. Consider a simple task of reorienting an object resting on a table with a two-fingered hand. To plan suitable contacts, it is necessary to reason about interaction forces. For example, a cube can be rotated by applying forces from the sides; however, if the object is a thin plate, the fingers must “press down” the object to create friction to achieve the same task. More importantly, as the fingers reach their respective kinematic limits, they need to break and re-establish contacts to rotate the object further. These two aspects, interaction forces and contact switches, have been the main challenges of object manipulation planning.

Over the past decade, trajectory optimization has been favored for multi-contact motion planning as this leads to efficient formulations to reason about interaction forces [1]–[4]. Yet, it remains unclear how the planning of contact modes should be efficiently incorporated, primarily due to its discrete nature that results in discontinuity in the dynamics at contact switch. To handle this discontinuity under the trajectory optimization framework, two main streams of methodologies have emerged:

- 1) the contact-invariant or contact-implicit approach enforces contact complementarity either as hard constraints [5], [6], penalty terms in the cost function [7]–[9], or with differentiable contact models [10], [11], and

- 2) the hybrid approach treats contact switches as discrete decisions and incorporates them in the continuous trajectory optimization problem [12]–[16].

In the context of robot manipulation, one representative work is Contact-Trajectory Optimization proposed in [13], where contact scheduling is modeled as binary decision variables and the non-convexity of the dynamics due to cross products is relaxed using McCormick envelopes [17]. The problem can then be formulated as Mixed-Integer Quadratic Programming (MIQP) [18]. However, the usage of the McCormick envelopes leads to a relaxed problem instead of the original one. As a result, applying the planned contact forces and locations may incur undesired torques. To date, the approach has only been demonstrated on 2D object manipulation with very short manipulation sequences, probably because its complexity grows exponentially with respect to the number of discrete variables due to the mixed-integer formulation.

Recent research has made progress towards speeding up Mixed-Integer Programming (MIP) [19] by leveraging machine learning techniques. For example, Nair et al. use neural networks to learn branch-and-bound heuristics and partial assignment for the discrete variables [20]; CoCo [21] finds feasible solution to MIP by learning to assign discrete variables and solving the resulted convex optimization problem. While such methodology greatly improves the solution speed at inference time, it assumes that the original MIP can be solved in a reasonable amount of time to construct the training set. If the original problem is prohibitive to solve, collecting a large dataset for this problem may not be practical without abundant computational resources or additional learning of a problem reduction [22].

In this work, we approach the problem from a different angle. Instead of modeling and solving it as a MIP, we formulate a tree search problem to find dynamically consistent contact sequences. Specifically, we use Monte Carlo Tree Search (MCTS) [23] to explore potential contact surfaces; then we formulate the underlying continuous optimization problem as a Biconvex Program [24] to allow efficient solution via the Alternating Direction Method of Multipliers (ADMM) [25] that has been adopted in online whole-body motion planning due to its guaranteed sublinear rate of convergence [26]. This leads to a formulation in which the discrete search space can be significantly reduced by introducing domain-specific heuristics for robot manipulation and the continuous problem can be solved efficiently without relaxation. More importantly, we show that learning-based MCTS trained on short-horizon tasks generalizes directly to long-horizon tasks. This removes the need for collecting data on large-scale problems, which can be time consuming. The

<sup>1</sup>Tandon School of Engineering, New York University, USA

This work was in part supported by the National Science Foundation (grants 1932187, 1925079, 2026479 and 2222815) and Meta Platforms, Inc.

contributions of the paper are as follows:

- 1) formulation of dynamically consistent contact planning for manipulation using learning-based MCTS,
- 2) efficient resolution of the underlying continuous optimization problem as a Biconvex Program with ADMM, and
- 3) extensive simulation including comparisons with MIQP approaches as well as real robot experiments to demonstrate the capabilities of our approach.

To our best knowledge, this is the first application of learning-based MCTS to dynamically consistent contact planning for manipulation.

## II. PROBLEM STATEMENT

### A. Inputs

We aim to solve an object manipulation task similar to the Contact-Trajectory Optimization problem proposed in [13] where the following quantities are given:

- 1) a rigid object with known geometry, friction coefficient  $\mu$ , mass  $m$ , moment of inertia  $\mathcal{I}$ , and  $N_\Omega$  pre-defined contact surfaces, each of which can be described as the convex span of its vertices  $v_{\Omega,i}$ ,
- 2) a trajectory with discretization step  $\Delta t$  of length  $T$  that consists of the desired object pose, velocity, and acceleration  $\xi = [q(t), \dot{q}(t), \ddot{q}(t)]_{t=1}^T$ , where  $q(t) = [p(t), R(t)] \in \text{SE}(3)$  consists of the position and orientation,  $\dot{q}(t) = [v(t), \omega(t)] \in \mathfrak{se}(3)$  consists of the linear and angular velocity, and  $\ddot{q}(t) = [\dot{v}(t), \dot{\omega}(t)]$  is the acceleration,
- 3) an environment with known geometry and friction coefficient  $\mu_e$ , and
- 4) a robot with  $N_c$  end-effectors that are able to make point contact with the object.

At the  $t$ -th time step, given the object motion and the object dynamics, we can compute the desired total force  $f_{\text{des}}(t)$  and torque  $\tau_{\text{des}}(t)$  to be applied to the object

$$\begin{bmatrix} f_{\text{des}}(t) \\ \tau_{\text{des}}(t) \end{bmatrix} = \begin{bmatrix} m(\dot{v}(t) + \omega(t) \times v(t) - g(t)) \\ \mathcal{I}\dot{\omega}(t) + \omega(t) \times \mathcal{I}\omega(t) \end{bmatrix}, \quad (1)$$

where all quantities, including the gravity term  $g(t)$  are expressed in the **object frame**.

In addition, as the geometry of the object and the environment as well as the object motion are known, we can obtain  $N_e(t)$  environment contact locations  $r_e(t)$  for  $e \in \{N_c+1, \dots, N_c+N_e(t)\}$  at each time step  $t$  by checking the collisions between the object and the environment, assuming uniform pressure distribution.

### B. Outputs

We aim to find the following:

- 1) the contact surface  $\Omega_c(t) \in \{0, 1, \dots, N_\Omega\}$ , the contact force  $f_c(t)$  and the contact location  $r_c(t)$  for each end-effector  $c$  of the robot;  $\Omega_c(t) = 0$  indicates that the  $c$ -th end-effector is not in contact, and
- 2) the environment contact force  $f_e(t)$

such that the forces and torques sum to the desired ones. Fig. 1 illustrates the outputs of our method for a double-finger manipulator pivoting a 2D square. Note that this is

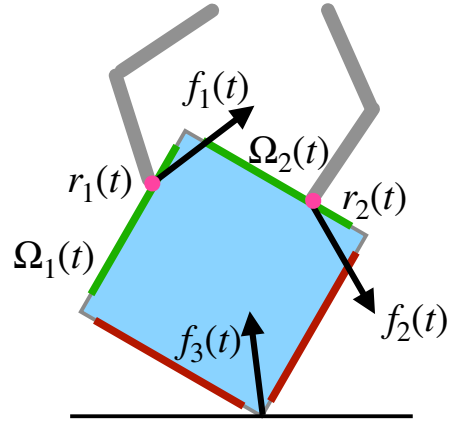


Fig. 1: Illustration of the outputs of the method.  $f_1(t), f_2(t)$  are the manipulator contact forces and  $r_1(t), r_2(t)$  are the respective contact locations (marked by the pink dots.)  $f_3(t)$  is the environment reaction force while the environment contact location is known beforehand. The thick lines on the object depict pre-defined contact surfaces and the green ones are the selected contact surfaces  $\Omega_1(t), \Omega_2(t)$  at time  $t$ .

only for illustrative purpose and our method is applicable to and tested on 3D objects and SE(3) poses as we will show in the experiments.

## III. METHOD

The problem described above is difficult even though the desired object motion is provided, as the solver needs to decide not only the contact force and location, but also the timing of contact switches. To address this challenge, we use learning-based MCTS to discover promising contact sequences and then evaluate their dynamical feasibility using an ADMM-based trajectory optimization algorithm. Fig. 2 provides an overview of our approach.

### A. Continuous Contact Optimization via ADMM

First, let us consider a simpler sub-problem where we already obtained a sequence of contact surfaces  $[\Omega_1(t), \dots, \Omega_{N_c}(t)]_{t=1}^T$  for each end-effector  $c$ . We can find the exact contact forces and locations by solving a non-convex continuous optimization problem. However, this non-convex problem, as we will show, can be formulated as a Biconvex Program and solved efficiently with ADMM. In contrast to [13], our formulation does not require piecewise linear approximation of the cross product and solves the exact original non-convex problem instead of a relaxed one. The optimization problem can be described by the following constraints and cost function:

1) *Dynamics*: The contact forces and torques must sum to the desired ones

$$\sum_{c=1}^{N_c} f_c(t) + \sum_{e=N_c+1}^{N_c+N_e(t)} f_e(t) = f_{\text{des}}(t) \quad (2a)$$

$$\sum_{c=1}^{N_c} r_c(t) \times f_c(t) + \sum_{e=N_c+1}^{N_c+N_e(t)} r_e(t) \times f_e(t) = \tau_{\text{des}}(t). \quad (2b)$$

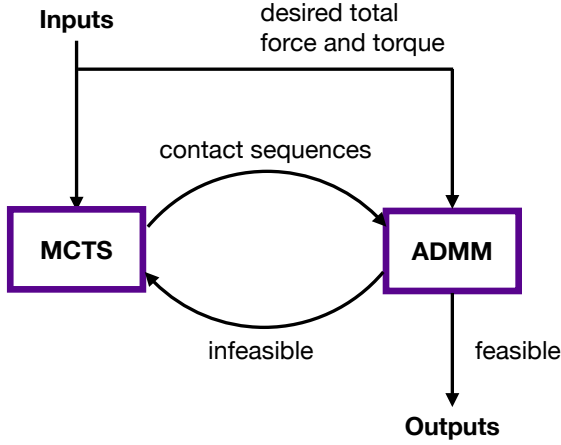


Fig. 2: An overview of the proposed method. First, candidate contact sequences are proposed by MCTS. Then, they are evaluated by an ADMM-based trajectory optimization algorithm to find dynamically feasible contact forces and locations to realize the desired object motion. At inference time, this repeats until the first feasible solution is found; at training time, we let the algorithm discover multiple solutions and collect both feasible and infeasible contact sequences to construct a diverse training set.

2) *Contact location*: The contact location must be inside the given contact surface  $\Omega_c(t)$  for  $\Omega_c(t) \neq 0$

$$\forall c \in \{c | \Omega_c(t) \neq 0\}, \sum_{i=1}^{N_{v,\Omega_c(t)}} \alpha_{c,i}(t) v_{\Omega_c(t),i} = r_c(t), \quad (3a)$$

$$\sum_{i=1}^{N_{v,\Omega_c(t)}} \alpha_{c,i}(t) = 1, \quad (3b)$$

where  $v_{\Omega_c(t),i}$  is the  $i$ -th vertex of the contact surface  $\Omega_c(t)$  for the end-effector  $c$  and  $\alpha_{c,i}(t) \geq 0$  is the corresponding coefficient of the convex span;  $N_{v,\Omega_c(t)}$  denotes the number of vertices.

3) *Contact force*: If the  $c$ -th end-effector is not in contact with any contact surface, hence  $\Omega_c(t) = 0$ , the contact force is set to zero

$$\forall c \in \{c | \Omega_c(t) = 0\}, f_c(t) = 0. \quad (4)$$

Note that this is not a complementarity constraint as  $\Omega_c(t)$  is already given by MCTS.

4) *Sticking contact*: To prevent the end-effector from sliding on the object, we impose that if the end-effector is in contact at time step  $t$ , it must remain sticking at time step  $t + 1$

$$\forall c \in \{c | \Omega_c(t) \neq 0\}, r_c(t + 1) = r_c(t). \quad (5)$$

Note that this constraint also implies that the end-effector cannot change its contact surface in one step; it has to break the current contact before switching to a new surface.

5) *Coulomb friction*: We assume all contact forces satisfy the Coulomb friction model. Hence, the sticking contact should stay inside the linearized friction cone of the given contact surface. The sliding contact (only environment contacts can slide) should satisfy  $f_e^\parallel(t) = -\mu_e \|f_e^\perp(t)\| \hat{r}_e^\parallel(t)$ , where  $f_e^\parallel(t)$  is the tangential force,  $f_e^\perp(t)$  the normal force, and  $\hat{r}_e^\parallel(t)$  the unit direction of the contact point velocity projected onto the contact surface, which can be obtained from the given object motion. For notational simplicity, we denote these constraints by the respective feasible set, hence

$$f_c(t) \in \mathbb{F}_c(t), f_e(t) \in \mathbb{F}_e(t) \quad (6)$$

6) *Cost function*: Finally, we minimize a quadratic objective function that avoids applying large forces at the boundary of the contact surface

$$J = \sum_{t=1}^T \sum_{c=1}^{N_c} \|f_c(t)\|^2 + \|r_c(t)\|^2 \quad (7)$$

7) *Biconvex Decomposition*: The optimization problem described above has an interesting feature that the only non-convex constraint (2b) due to the cross product  $r_c(t) \times f_c(t)$  is in fact biconvex. That is, when  $r_c(t)$  is fixed, this constraint is convex with respect to  $f_c(t)$ ; when  $f_c(t)$  is fixed, this constraint is convex with respect to  $r_c(t)$ . Note that in both cases, the environment contact location  $r_e(t)$  is known and thus the cross product  $r_e(t) \times f_e(t)$  does not pose any non-convexity. In fact, when we group the decision variables into two sets  $x = [r_c(t), \alpha_c(t)]_{t=1}^T$  and  $z = [f_c(t), f_e(t)]_{t=1}^T$ , we can re-write the original problem into the standard ADMM form with the constraint

$$G(x, z) = \sum_{c=1}^{N_c} r_c(t) \times f_c(t) + \sum_{e=N_c+1}^{N_c+N_e(t)} r_e(t) \times f_e(t) - \tau_{\text{des}}(t) = 0,$$

and the iteration

$$x^{k+1} = \arg \min_x \sum_{t=1}^T \sum_{c=1}^{N_c} \|r_c(t)\|^2 + \frac{\rho}{2} \|G(x, z^k) + y^k\| \quad (8a)$$

s.t. (3), (5)

$$z^{k+1} = \arg \min_z \sum_{t=1}^T \sum_{c=1}^{N_c} \|f_c(t)\|^2 + \frac{\rho}{2} \|G(x^{k+1}, z) + y^k\| \quad (8b)$$

s.t. (2a), (4), (6)

$$y^{k+1} = y^k + G(x^{k+1}, z^{k+1}), \quad (8c)$$

where  $y$  is the scaled dual variable and  $\rho > 0$  is a penalty parameter that is tuned to  $2 \times 10^6$  in the experiments. The ADMM iteration initializes with  $\alpha_{c,i}^0(t) = 1/N_{v,\Omega_c(t)}$ ,  $r_c^0(t) = \sum_{i=1}^{N_{v,\Omega_c(t)}} \alpha_{c,i}^0(t) v_{\Omega_c(t),i}$ ,  $f_c^0(t) = 0$ ,  $f_e^0(t) = 0$  and  $y^0 = 0$ . Note that both the optimization problems (8a) and (8b) are just Quadratic Programs (QPs) [27] which are simple to solve. Throughout our experiments, we only run one ADMM iteration as we observe satisfactory convergence in this setting. Hence, solving this non-convex optimization problem only requires solving two QPs.

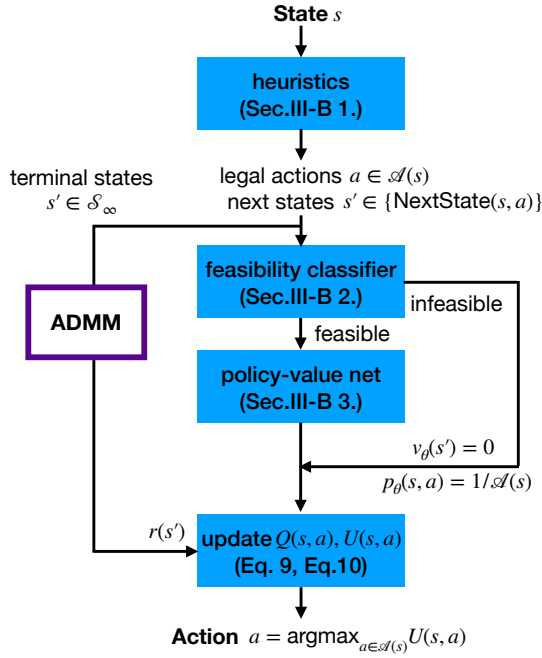


Fig. 3: Action selection in the MCTS search process. The actions are evaluated by both domain-specific heuristics and trained neural networks to enable efficient search. As the search reaches a terminal state, the reward is computed by ADMM. For readability, recursions are omitted.

### B. Discrete Contact Planning via MCTS

Now that we have shown the contact locations and forces can be found efficiently if the contact surfaces are known, we focus on the missing piece: finding these contact surfaces. To do this, we adopt MCTS that was behind many recent advances in reinforcement learning [23], [28] as well as in robotic applications such as gait discovery [29], task and motion planning [30], and object rearrangement [31], [32]. To solve the contact sequence discovery problem via MCTS, we first model it as a Markov-decision Process (MDP) with states  $s \in \mathcal{S}$  and actions  $a \in \mathcal{A}(s)$ , where  $\mathcal{S}$  is the set of states and  $\mathcal{A}(s)$  is the set of legal actions at the state  $s$ . At time step  $k$ , the action  $a_k = [\Omega_1(k), \dots, \Omega_{N_c}(k)]$  represents which object surface each end-effector needs to be in contact with (or not) and the state  $s_k$  contains the current desired object pose  $q(k)$  and the history of the object surfaces that were in contact with each end-effector  $[\Omega_1(t), \dots, \Omega_{N_c}(t)]_{t=1}^k$ . Note that the state transition  $s_{k+1} = \text{NEXTSTATE}(s_k, a_k)$  in this MDP is deterministic as the desired pose is known a priori and thus fixed. With this formulation, we modify the MCTS algorithm by utilizing both domain-specific heuristics and neural networks trained from past experience to improve search efficiency. As the search reaches the end, an optimization problem is constructed and solved by ADMM to provide a terminal reward. Fig. 3 previews these modifications that will be elaborated in the remainder of this section.

Let us first recall the standard learning-based MCTS algorithm as summarized in Algorithm 1. It constructs a search tree  $\mathcal{T} = (\mathcal{V}, \mathcal{E})$  where the set of vertices  $\mathcal{V}$  contains

the visited states and the set of edges  $\mathcal{E}$  contains the visited transitions  $(s \xrightarrow{a} s')$ . For each transition, it maintains the state-action value  $Q(s, a)$  that estimates future rewards to be accumulated and the number of occurrences  $N(s, a)$  of this state-action pair during the search. To update the state-action value  $Q(s, a)$ , learning-based MCTS uses a policy-value network (parameterized by  $\theta$ ) to provide a value estimate  $v_\theta(s')$  for the possible next states  $s' \in \{\text{NEXTSTATE}(s, a) | a \in \mathcal{A}(s)\}$

$$Q(s, a) \leftarrow \frac{N(s, a)Q(s, a) + v_\theta(s')}{N(s, a) + 1}. \quad (9)$$

The same network also outputs an action probability  $p_\theta(s, a)$  to calculate a heuristic score for the state-action pair

$$U(s, a) = Q(s, a) + \gamma p_\theta(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}, \quad (10)$$

where  $\gamma > 0$  is a hyper-parameter controlling the degree of exploration; in our experiments, it is manually tuned to 0.1. This score is used by MCTS to select promising actions  $a = \arg \max_{a \in \mathcal{A}(s)} U(s, a)$  while balancing exploration with exploitation. Once the search reaches a terminal state  $s \in \mathcal{S}_\infty$ , the contact surfaces found by MCTS are used to construct the optimization problem described in Sec. III-A. Its solution will then be evaluated to return a reward  $r$  to update the state-action value and guide future search.

To calculate the reward, we integrate the solution to obtain the final object pose  $\hat{q}(T) = [\hat{p}(T), \hat{R}(T)]$  with the semi-implicit Euler method. We then compare it with the desired final pose  $q(T) = [p(T), R(T)]$  to compute a weighted distance

$$D(q, \hat{q}) = \|p - \hat{p}\| + \beta \left\| \log(\hat{R}^T R) \right\|, \quad (11)$$

where  $\beta > 0$  scales the angular distance; in the experiments, it is set to 0.1. Note that this distance does not always equal zero as we terminate ADMM only one iteration. The weighted distance within a threshold  $D \leq D_{\text{th}}$  is then normalized to  $[0, 1]$  to obtain the reward. In our experiments, the threshold  $D_{\text{th}}$  is set to 0.03 which allows a maximal object position error of 3cm or an orientation error of 0.3rad. Note that we set a relatively high threshold in order to collect diverse training data; after training, the MCTS almost always discovers solutions with near zero error as will be shown in the experiments. If the contact sequence found by MCTS does not lead to a dynamically feasible solution or has a higher error above the threshold, the reward will be set to zero. Hence, the EVALUATE( $s$ ) function in MCTS entails the ADMM iteration (8) and computing the reward.

At inference time, the procedure  $\text{MCTS}(s; \theta)$  is repeated until the first feasible solution is found; at training time, we let it discover multiple solutions and collect both feasible and infeasible contact sequences to construct a diverse training set  $\mathcal{D} = \{(\bar{v}(s), \bar{p}(s, a)) | s \in \mathcal{V}\}$  for all visited states, where  $\bar{p}(s, a) = N(s, a) / \sum_b N(s, b)$  is the empirical action probability and  $\bar{v}(s) = \sum_{a \in \mathcal{A}(s)} \bar{p}(s, a) Q(s, a)$  is the expected state value. The network is then updated using the sum of a

---

**Algorithm 1** Learning-Based Monte Carlo Tree Search

---

```

1: procedure MCTS( $s; \theta$ )
2:   if  $s \in \mathcal{S}_\infty$  then
3:      $r \leftarrow \text{EVALUATE}(s)$ 
4:     return  $r$ 
5:   else if  $s \notin \mathcal{V}$  then
6:      $\mathcal{V} \leftarrow \mathcal{V} \cup \{s\}$ 
7:     for  $a \in \mathcal{A}(s)$  do
8:        $Q(s, a) \leftarrow 0$ 
9:        $N(s, a) \leftarrow 0$ 
10:    end for
11:    return  $v_\theta(s)$ 
12:   else
13:      $a \leftarrow \arg \max_{a \in \mathcal{A}(s)} U(s, a)$ 
14:      $s' \leftarrow \text{NEXTSTATE}(s, a)$ 
15:      $v \leftarrow \text{MCTS}(s'; \theta)$ 
16:      $Q(s, a) \leftarrow \frac{N(s, a)Q(s, a) + v}{N(s, a) + 1}$ 
17:      $N(s, a) \leftarrow N(s, a) + 1$ 
18:      $\mathcal{E} \leftarrow \mathcal{E} \cup \{(s \xrightarrow{a} s')\}$ 
19:     return  $v$ 
20:   end if
21: end procedure

```

---

mean-square loss  $l_v$  for the value head and a cross-entropy loss  $l_p$  for the policy head

$$l(\mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{s \in \mathcal{V}} (l_v(s) + l_p(s)), \quad (12)$$

where

$$l_v(s) = (v_\theta(s) - \bar{v}(s))^2 \quad (13a)$$

$$l_p(s) = - \sum_{a \in \mathcal{A}(s)} \bar{p}(s, a) \log p_\theta(s, a). \quad (13b)$$

*1) Assumptions and Heuristics:* One major advantage of using MCTS over MIQP is that it is straightforward to incorporate domain-specific assumptions and heuristics to reduce the search space. In this work, we apply the following assumptions and heuristics:

- **Contact surface:** Each contact surface can only be touched by at most one end-effector and each end-effector can touch at most one contact surface.
- **Persistent contact:** While the downstream continuous optimization problem may have a small discretization step, for example  $\Delta t = 0.1$  s, most manipulation tasks do not require decisions of contact switch at such a high resolution. Thus, we assume that an end-effector must remain on the same surface for  $d$  time steps, which means a trajectory of length  $T$  can admit at most  $T/d - 1$  contact switches.
- **Kinematic feasibility:** For each end-effector  $c$ , a contact surface will only be considered if inverse kinematics can find a robot configuration that reaches the center of this surface within an error threshold of 2 cm and does not result in any undesired collision (e.g. between non-end-effector links and the object). Note that this cannot be imposed in a MIQP formulation as it introduces nonlinear constraints.
- **Contact switch:** We allow at most one end-effector to make or break the contact at each time step. Moreover, the end-effector can only break the contact if the desired object velocity and acceleration is zero.

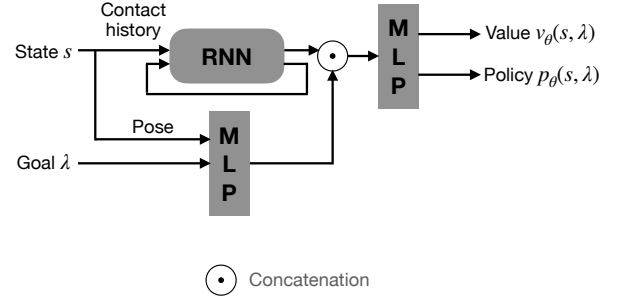


Fig. 4: Schematic diagram of the policy-value network architecture. Activation functions and regularization layers such as Dropout and BatchNorm are omitted.

*2) Feasibility classifier:* One key difference between our task and generic game-play is that our dataset is highly imbalanced—many contact sequences explored by the MCTS are dynamically infeasible, resulting in zero reward. Directly training on such a dataset leads to underestimation of the value function. Instead, we first train on the dataset  $\mathcal{D}$  a binary classifier  $C_\phi(s)$  with logistic regression where dynamically feasible samples are given more weights. At inference time, a state is only fed into the policy-value network if the classifier labels it as dynamically feasible; otherwise, we output zero value  $v_\theta(s) = 0$  and uniformly distributed action probability  $p_\theta(s, a) = 1/|\mathcal{A}(s)|$ . This feasibility classifier screens out dynamically infeasible contact sequences before the MCTS completes the search, thus greatly improves search efficiency.

*3) Goal-conditioned policy-value network:* Note that each MCTS instance only searches for the contact sequence for a given object motion  $\xi$ , thus the rewards are motion-specific. If we were to learn from the data collected for this object motion only, it is unlikely that the network would generalize to other motions. Thus, we generate multiple object motions in the training phase and additionally input an intermediate goal variable to the network. It is defined as the difference between the current desired pose and the one  $h$  steps in the future  $\lambda(t) = q(t+h) \ominus q(t)$ , where  $\ominus$  denotes subtraction in SE(3). Fig. 4 depicts the policy-value network architecture: it takes as inputs the state  $s$  and the goal  $\lambda$ , and outputs the goal-conditioned value  $v_\theta(s, \lambda)$  and action probabilities  $p_\theta(s, \lambda)$ . Since the sequence of contact surfaces has varying lengths, we use a Recurrent Neural Network (RNN) to encode this information and concatenate it with the pose and the goal processed by a Multilayer Perceptron (MLP). Due to the usage of the feasibility classifier mentioned above, we only train our policy-value network on a subset  $\mathcal{D}_+ \subseteq \mathcal{D}$  with positive samples  $\mathcal{V}_+ = \{s \in \mathcal{V} | \bar{v}(s) > 0\}$  to avoid underestimating the value function.

#### IV. EXPERIMENTS

We conduct experiments in simulation and on real hardware to show that our method 1) is capable of finding high quality dynamically feasible solutions much faster than a MIQP baseline. 2) scales to long-horizon tasks even when trained only on data collected from shorter-horizon tasks.



### A. Experiment Setup

1) *Tasks*: Throughout all experiments, we consider a manipulator composed of two modular robot fingers similar to the ones used in [33] and a  $10\text{ cm} \times 10\text{ cm} \times 10\text{ cm}$  cube with mass  $m = 0.5\text{ kg}$  on an infinitely large plane. The cube and the plane have the same friction coefficient  $\mu = \mu_e = 0.8$ . We consider one contact surface for each face of the cube except for the bottom one; each contact surface is a  $8\text{ cm} \times 8\text{ cm}$  square to avoid contact locations at the corner. The object motion trajectory is generated with spline interpolation in SE(3) between the initial object pose and a desired pose parameterized as the following primitives: 1) **sliding** ( $S$ ) on the  $xy$ -plane by  $-10\text{ cm}$  to  $10\text{ cm}$  2) **sliding with curvature** ( $SC$ ) on the  $xy$ -plane by  $-5\text{ cm}$  to  $5\text{ cm}$  with a rotation about the  $z$ -axis by  $-45^\circ$  to  $45^\circ$  3) **rotating** ( $R$ ) about the  $z$ -axis by  $-90^\circ$  to  $90^\circ$  4) **lifting** ( $L$ ) along the  $z$ -axis by  $0\text{ cm}$  to  $10\text{ cm}$ , and 5) **pivoting** ( $P$ ) about the  $y$ -axis by  $0^\circ$  to  $45^\circ$ . The  $xy$ -axes span the plane that the object is placed on and the  $z$ -axis points to the opposite gravity direction. Finally, the initial object position is randomly sampled on the  $xy$ -plane within a  $[-5\text{ cm}, 5\text{ cm}]^2$  area centered at the origin and the initial orientation about the  $z$ -axis by  $-90^\circ$  to  $90^\circ$ .

2) *Baselines*: We compare our method (**MCTS**) with two baselines:

- the **MIQP** baseline is implemented following [13]. We did not use the authors' open-source implementation as it was only implemented for 2D objects. But the same formulation can be directly extended to 3D. To facilitate a fair comparison, we also implemented all heuristics described in Sec III-B.1 except the kinematic feasibility check. For the McCormick envelope relaxation of the cross product, we partition the contact location into 4 intervals and the contact force into 2 intervals. In all experiments, we terminate the MIQP solver at the first feasible solution instead of waiting for the global optimum which may take extremely long time. In addition, we implement the constraint (2a) as a penalty term in the cost function. This accelerates the MIQP solver significantly from our observation during the experiments. We note that our implementation has comparable computation time as reported in [13].
- the **MCTS<sup>U</sup>** baseline represents an untrained model and constantly outputs zero values  $v_\theta(s) = 0$  and uniform action probability  $p_\theta(s, a) = 1/|\mathcal{A}(s)|$ .

3) *Software and hardware*: We conduct all experiments on a single GeForce RTX 2080 TI GPU and an Intel Xeon CPU at 3.7 GHz using Python and PyTorch. We model and solve the QPs with CVXPY [34] and OSQP [35] and use Gurobi [36] for MIQP. All source code including the baseline can be found at <https://hzhu.io/contact-mcts>.

### B. Evaluation metrics

We examine two metrics to evaluate the effectiveness and efficiency of our method: 1) the **force and torque error** between the desired and the solution. The error is averaged over the horizon  $T$  and scaled by the object mass and inertia

respectively. The smaller this error is, the better the solution tracks the desired object motion. 2) The **computation time** needed to find the first dynamically feasible solution.

### C. Training procedure

We generate 300 object motion trajectories, each comprising two primitives with randomly sampled parameters. In particular, 200 trajectories are composed of two  $SC$  primitives; 50 trajectories of one  $SC$  and one  $L$ ; 50 trajectories of one  $SC$  and one  $P$ . For the  $i$ -th trajectory, we let an untrained MCTS run until it evaluates 200 candidate contact sequences; then we construct the dataset  $\mathcal{D} = \mathcal{D} \cup \{(\bar{v}(s), \bar{p}(s)) | s \in \mathcal{V}_i\}$  where  $\mathcal{V}_i$  contains all the states the MCTS visited for the  $i$ -th trajectory. The policy-value network and the value classifier are then trained via Adam [37] for 300 epochs.

### D. Single motion primitives

In this experiment, we consider object motion trajectories that consist of one single primitive. Each primitive has a desired pose uniformly randomly sampled from its respective parameter range described in Sec. IV-A.1. Each trajectory consists of  $T = 10$  time steps with step size  $\Delta t = 0.1\text{ s}$ ; each contact persists as well at least  $0.1\text{ s}$ , hence a trajectory can admit at most 9 contact switches. We run 50 trials for each primitive to collect the performance statistics.

Table I shows that our method is capable of finding dynamically feasible solutions consistently faster than the MIQP baseline thanks to the MCTS formulation. Especially for primitives that require non-zero torques ( $R$ ,  $SC$ ,  $P$ ), the MIQP baseline is not only an order of magnitude slower, but also produces solutions with large errors. We note that the force error can be reduced by letting the MIQP solver explore more feasible solutions, while the torque error remains high nonetheless. This might be due to its usage of the McCormick envelopes to approximate cross products, which not only introduces approximation error but also adds additional discrete variables. In contrast, thanks to the ADMM formulation, our method solves the original problem instead of a relaxed one and has thus near-zero average force and torque error.

We also note that while the solutions found by the MIQP baseline are dynamically feasible, they are not necessarily kinematically feasible or collision-free since these conditions cannot be incorporated as linear constraints. While it is possible to collect multiple dynamically feasible solutions and pick the kinematically feasible one from them, it may further increase the computation time.

### E. Long-horizon tasks

In the previous experiments, we have shown the effectiveness of our method for short motion primitives. Let us now consider tasks that last a longer period of time. First, we extend the primitive to  $T = 30$  time steps by stipulating each contact persists for  $d = 3$  steps. Note that there are still at most 9 contact switches for a single primitive. However, such extended primitives may be useful for tasks that require longer execution time but not necessarily more contact switches; for instance, sliding a cube for a long distance

TABLE I: Task performance for object motion primitives. Values  $\leq 0.005$  are rounded to zero.

	Method	Time [s]		Error [N, N · m]	
		Mean	Worst	Mean	Worst
<i>S</i>	MIQP	0.65	0.79	0.66, <b>0.00</b>	2.90, <b>0.00</b>
	MCTS	<b>0.10</b>	<b>0.18</b>	<b>0.00, 0.00</b>	<b>0.00, 0.00</b>
	MCTS <sup>U</sup>	0.24	1.23	<b>0.00</b> , 0.03	0.06, 1.14
<i>L</i>	MIQP	0.25	0.51	6.29, <b>0.00</b>	6.87, <b>0.00</b>
	MCTS	<b>0.15</b>	<b>0.23</b>	<b>0.24</b> , 0.05	<b>0.86</b> , 0.18
	MCTS <sup>U</sup>	0.53	2.23	0.53, 0.11	0.88, 0.35
<i>R</i>	MIQP	4.83	29.46	8.36, 16.64	30.72, 45.57
	MCTS	<b>0.12</b>	<b>0.27</b>	<b>0.00, 0.00</b>	<b>0.00, 0.00</b>
	MCTS <sup>U</sup>	0.41	1.22	<b>0.00, 0.00</b>	<b>0.00, 0.00</b>
<i>SC</i>	MIQP	2.19	4.41	11.73, 22.39	49.61, 88.27
	MCTS	<b>0.11</b>	<b>0.24</b>	<b>0.00, 0.00</b>	<b>0.00, 0.00</b>
	MCTS <sup>U</sup>	0.20	0.81	<b>0.00, 0.00</b>	0.03, 0.07
<i>P</i>	MIQP	6.69	50.41	7.65, 15.31	26.85, 53.65
	MCTS	<b>0.15</b>	<b>0.34</b>	<b>0.01, 1.23</b>	<b>0.23, 14.01</b>
	MCTS <sup>U</sup>	0.17	0.45	<b>0.01</b> , 1.46	0.33, 19.55

or rotating it very slowly. In the following experiments, we concatenate such extended primitives to form a even longer trajectory. In particular, we consider the primitive *SC* as it represents typical planar repositioning/reorienting tasks.

Table II reports the performance metrics for each method. A task is considered failed if no dynamically feasible solution is found within 60 s. We can immediately see that the trained MCTS efficiently solves all the tasks regardless of the trajectory length, while the MIQP baseline and the untrained MCTS struggle in long-horizon tasks (the errors decrease because they are computed only on successful trials). Indeed, the MIQP baseline cannot solve any tasks containing more than two primitives in 60 s. Interestingly, even for the task with a single extended primitive, where the number of possible contact switches does not change compared to the previous tasks in Sec. IV-D, the MIQP baseline still need significantly more time to find a feasible solution. This could again be attributed to the McCormick envelopes as they add additional discrete variables to each time step regardless of the underlying number of contact switches.

Finally, we highlight that the MCTS training dataset only contains object motion trajectories consisting of at most two primitives. However, Table II shows that our method is able to efficiently solve the longer-horizon tasks without being explicitly trained on them as our MCTS formulation exploits the intrinsic compositionality of the task by learning a goal-conditioned policy-value network. Hence, we do not need to collect training data on large-scale, time-consuming problems as opposed to the learning-based MIP method proposed in [20].

#### F. Executing the contact plan

To validate the contact plans found by our method, we execute them in an open-loop fashion with a simple impedance controller in the PyBullet simulator [38] and on a real robot

$$\tau = J^T (K(r_c^w - r^w) + D(\dot{r}_c^w - \dot{r}^w) + f_c^w), \quad (14)$$

where  $J$  is the end-effector Jacobian;  $K, D$  are manually tuned gain matrices;  $r^w, \dot{r}^w$  are the position and velocity of

TABLE II: Task performance for object motions composed of *SC* primitives. Errors are computed only on successful trials. Values  $\leq 0.005$  are rounded to zero.

# SC	Method	Success rate	Time [s]		Error [N, N · m]	
			Mean	Worst	Mean	Worst
1	MIQP	94%	10.15	60.00	3.40, 11.72	19.93, 41.68
	MCTS	<b>100%</b>	<b>0.21</b>	<b>0.41</b>	<b>0.00, 0.00</b>	<b>0.00, 0.00</b>
	MCTS <sup>U</sup>	<b>100%</b>	0.91	3.67	<b>0.00, 0.00</b>	0.03, 0.07
2	MIQP	42%	40.93	60.00	4.96, 4.38	16.61, 22.54
	MCTS	<b>100%</b>	<b>0.47</b>	<b>1.56</b>	<b>0.00, 0.00</b>	<b>0.01, 0.03</b>
	MCTS <sup>U</sup>	<b>100%</b>	3.08	12.84	<b>0.00, 0.00</b>	0.03, 0.07
3	MIQP	0%	—	—	—	—
	MCTS	<b>100%</b>	<b>1.35</b>	<b>8.84</b>	<b>0.00, 0.00</b>	<b>0.01, 0.03</b>
	MCTS <sup>U</sup>	90%	20.87	60.00	<b>0.00, 0.00</b>	0.01, 0.04

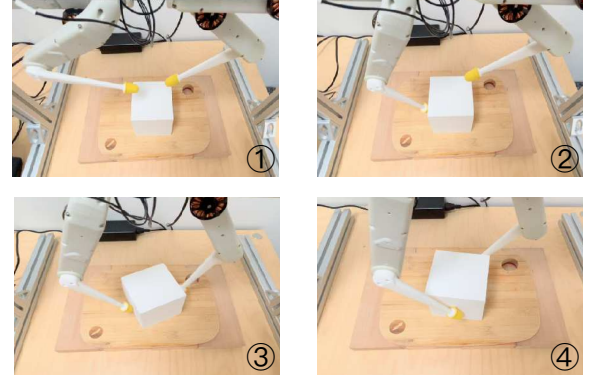


Fig. 5: Exemplary execution of rotating the cube by 90°. For a video compilation of various tasks, please refer to the supplemental materials.

the end-effector and  $f_c^w, r_c^w, \dot{r}_c^w$  are the planned contact force, location and velocity, all expressed in the world frame.

Fig. 5 shows an example of the contact plan execution of rotating a cube by 90°. The robot is able to move the object towards the target pose if the object is placed at the desired initial position. We note that the same impedance is used for all tasks with different primitives. This would not be possible for us without applying the planned contact forces; for example, the fingers would drop the cube if they were to lift it with purely position commands and too low of a position gain. This shows the benefits and importance of planning accurate forces and torques. However, we do observe failure cases when the reality differs too much from the model, especially for the pivoting primitive *P* which is sensitive to the discrepancy between the modeled and actual friction. We show this in the supplementary materials and point out the necessity of incorporating sensory feedback for online re-planning, which we leave for future work.

#### V. CONCLUSION

In this work, we proposed a framework that combines data-driven MCTS and efficient non-convex optimization via ADMM to find dynamically feasible contact forces and locations to realize a given object motion. We show that the

capability of learning from data allows our method to achieve great scalability for long-horizon tasks even when the dataset only contains short-horizon data.

The most limited aspect of our approach is that the object motion must be provided. True dexterity requires automatic discovery of object motion together with the contact plan. One potential way to achieve this is to jointly enumerate manipulator and environment contacts [15]. Another problem is that we represent contact surfaces as integers. While this is natural with the MCTS formulation, it makes the learned networks object-specific. To address this issue, it might be helpful to map the integer-valued surfaces to its geometric center before feeding them into the neural networks. Finally, our approach assumes perfect knowledge of the object and the environment, which is not possible in the real world. Thus, it is necessary to explore ways of integrating perception, for example, as done in [39].

## REFERENCES

- [1] A. Escande and A. Kheddar, "Contact planning for acyclic motion with tasks constraints," *IEEE/RSJ International Conference on Intelligent Robots and Systems, 2009. IROS 2009.*, p. 435–440, Oct 2009.
- [2] Y.-C. Lin, L. Righetti, and D. Berenson, "Robust humanoid contact planning with learned zero- and one-step capturability prediction," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, 2020.
- [3] B. Ponton, M. Khadiv, A. Meduri, and L. Righetti, "Efficient multi-contact pattern generation with sequential convex approximations of the centroidal dynamics," *IEEE Transactions on Robotics*, vol. 37, no. 5, p. 1661–1679, 2021.
- [4] J. Carpentier, S. Tonneau, M. Naveau, O. Stasse, and N. Mansard, "A versatile and efficient pattern generator for generalized legged locomotion," in *IEEE-RAS International Conference on Robotics and Automation*, 2016.
- [5] D. Stewart and J. C. Trinkle, "An implicit time-stepping scheme for rigid body dynamics with coulomb friction," in *IEEE International Conference on Robotics and Automation*, 2000, pp. 162–169.
- [6] M. Posa, C. Cantu, and R. Tedrake, "A direct method for trajectory optimization of rigid bodies through contact," *The Int J of Robotics Research*, vol. 33, no. 1, pp. 69–81, 2014.
- [7] I. Mordatch, E. Todorov, and Z. Popović, "Discovery of complex behaviors through contact-invariant optimization," *ACM Transactions on Graphics (TOG)*, vol. 31, no. 4, pp. 1–8, 2012.
- [8] I. Mordatch, J. M. Wang, E. Todorov, and V. Koltun, "Animating human lower limbs using contact-invariant optimization," *ACM Transactions on Graphics (TOG)*, vol. 32, no. 6, pp. 1–8, 2013.
- [9] I. Mordatch, Z. Popović, and E. Todorov, "Contact-invariant optimization for hand manipulation," in *ACM SIGGRAPH/Eurographics symposium on computer animation*, 2012, pp. 137–144.
- [10] M. Neunert et al., "Whole-body nonlinear model predictive control through contacts for quadrupeds," *IEEE Robotics and Automation Letters*, vol. 3, no. 3, pp. 1458–1465, 2018.
- [11] T. Pang, H. Suh, L. Yang, and R. Tedrake, "Global planning for contact-rich manipulation via local smoothing of quasi-dynamic contact models," *arXiv preprint arXiv:2206.10787*, 2022.
- [12] M. Toussaint, "Logic-geometric programming: An optimization-based approach to combined task and motion planning," in *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.
- [13] B. Aceituno-Cabezas and A. Rodriguez, "A global quasi-dynamic model for contact-trajectory optimization," in *Robotics: Science and Systems (RSS)*, 2020.
- [14] N. Doshi, F. R. Hogan, and A. Rodriguez, "Hybrid differential dynamic programming for planar manipulation primitives," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 6759–6765.
- [15] X. Cheng, E. Huang, Y. Hou, and M. T. Mason, "Contact mode guided motion planning for quasidynamic dexterous manipulation in 3d," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 2730–2736.
- [16] C. Chen, P. Culbertson, M. Lepert, M. Schwager, and J. Bohg, "Trajectory optimization meets tree search for planning multi-contact dexterous manipulation," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2021, pp. 8262–8268.
- [17] G. P. McCormick, "Computability of global solutions to factorable nonconvex programs: Part i—convex underestimating problems," *Mathematical programming*, vol. 10, no. 1, pp. 147–175, 1976.
- [18] R. Lazimy, "Mixed-integer quadratic programming," *Mathematical Programming*, vol. 22, pp. 332–349, 1982.
- [19] C. A. Floudas, *Nonlinear and mixed-integer optimization: fundamentals and applications*. Oxford University Press, 1995.
- [20] V. Nair, S. Bartunov, F. Gimeno, I. von Glehn, P. Lichocki, I. Lobov, B. O'Donoghue, N. Sonnerat, C. Tjandraatmadja, P. Wang, et al., "Solving mixed integer programs using neural networks," *arXiv preprint arXiv:2012.13349*, 2020.
- [21] A. Cauligi et al., "Coco: Online mixed-integer control via supervised learning," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 1447–1454, 2021.
- [22] X. Lin, G. I. Fernandez, and D. W. Hong, "Reduce: Reformulation of mixed integer programs using data from unsupervised clusters for learning efficient strategies," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 4459–4465.
- [23] D. Silver et al., "Mastering the game of go without human knowledge," *Nature*, vol. 550, no. 7676, pp. 354–359, 2017.
- [24] J. Gorski, F. Pfeuffer, and K. Klamroth, "Biconvex sets and optimization with biconvex functions: a survey and extensions," *Mathematical methods of operations research*, vol. 66, no. 3, pp. 373–407, 2007.
- [25] S. Boyd et al., "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends® in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [26] A. Meduri, P. Shah, J. Viereck, M. Khadiv, I. Havoutis, and L. Righetti, "Biconmp: A nonlinear model predictive control framework for whole body motion planning," *IEEE Transactions on Robotics*, 2023.
- [27] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [28] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, et al., "Mastering atari, go, chess and shogi by planning with a learned model," *Nature*, vol. 588, no. 7839, pp. 604–609, 2020.
- [29] L. Amatucci, J.-H. Kim, J. Hwangbo, and H.-W. Park, "Monte carlo tree search gait planner for non-gaited legged system control," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 4701–4707.
- [30] T. Ren, G. Chaltatzaki, and J. Peters, "Extended tree search for robot task and motion planning," *arXiv preprint arXiv:2103.05456*, 2021.
- [31] Y. Labbé, S. Zagoruyko, I. Kalevatykh, I. Laptev, J. Carpentier, M. Aubry, and J. Sivic, "Monte-carlo tree search for efficient visually guided rearrangement planning," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3715–3722, 2020.
- [32] H. Song, J. A. Haustein, W. Yuan, K. Hang, M. Y. Wang, D. Kragic, and J. A. Stork, "Multi-object rearrangement with monte carlo tree search: A case study on planar nonprehensile sorting," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 9433–9440.
- [33] M. Wüthrich et al., "Trifinger: An open-source robot for learning dexterity," *arXiv preprint arXiv:2008.03596*, 2020.
- [34] S. Diamond and S. Boyd, "CVXPY: A Python-embedded modeling language for convex optimization," *Journal of Machine Learning Research*, vol. 17, no. 83, pp. 1–5, 2016.
- [35] B. Stellato, G. Banjac, P. Goulart, A. Bemporad, and S. Boyd, "OSQP: an operator splitting solver for quadratic programs," *Mathematical Programming Computation*, vol. 12, no. 4, pp. 637–672, 2020.
- [36] Gurobi Optimization, LLC, "Gurobi Optimizer Reference Manual," 2022. [Online]. Available: <https://www.gurobi.com>
- [37] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [38] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," <http://pybullet.org>.
- [39] D. Driess, J.-S. Ha, and M. Toussaint, "Learning to solve sequential physical reasoning problems from a scene image," *The Int. J of Robotics Research*, vol. 40, no. 12–14, pp. 1435–1466, 2021.