

# Enhancing Task Performance of Learned Simplified Models via Reinforcement Learning

Hien Bui<sup>1</sup> and Michael Posa<sup>1</sup>

**Abstract**—In contact-rich tasks, the hybrid, multi-modal nature of contact dynamics poses great challenges in model representation, planning, and control. Recent efforts have attempted to address these challenges via data-driven methods, learning dynamical models in combination with model predictive control. Those methods, while effective, rely solely on minimizing forward prediction errors to hope for better task performance with MPC controllers. This weak correlation can result in data inefficiency as well as limitations to overall performance. In response, we propose a novel strategy: using a policy gradient algorithm to find a simplified dynamics model that explicitly maximizes task performance. Specifically, we parameterize the stochastic policy as the perturbed output of the MPC controller, thus, the learned model representation can directly associate with the policy or task performance. We apply the proposed method to contact-rich tasks where a three-fingered robotic hand manipulates previously unknown objects. Our method significantly enhances task success rate by up to 15% in manipulating diverse objects compared to the existing method while sustaining data efficiency. Our method can solve some tasks with success rates of 70% or higher using under 30 minutes of data. All videos and codes are available at <https://sites.google.com/view/lcs-rl>.

## I. INTRODUCTION

In many robotics tasks, such as dexterous manipulation and locomotion, robots frequently need to make and break contact with the environment. Yet, finding explicit models and policies that can exploit the hybrid complex interaction of the robot with its environment to solve the tasks remains a challenge. Some works in model-based control [1]–[4] have attempted to explicitly identify contact modes and plan the contact sequences. However, these approaches face scalability challenges as the number of contact modes increases.

Recent data-driven methods have made significant advances in tackling that scalability issue, broadly offering two primary strategies. Modern model-free reinforcement learning (RL) directly parameterizes control policies with deep neural networks, then iteratively improves the policies through large-scale trial and error [5]–[7]. However, because of their data inefficiency, carrying out experiments on real robotic systems is always resource-intensive and time-consuming. In contrast, a large portion of model-based RL [8]–[14] leverages the expressiveness power of deep neural networks to learn intricate dynamic models. The learned models are subsequently employed for trajectory planning/predictive control, typically via random shooting techniques. Despite improving data efficiency compared

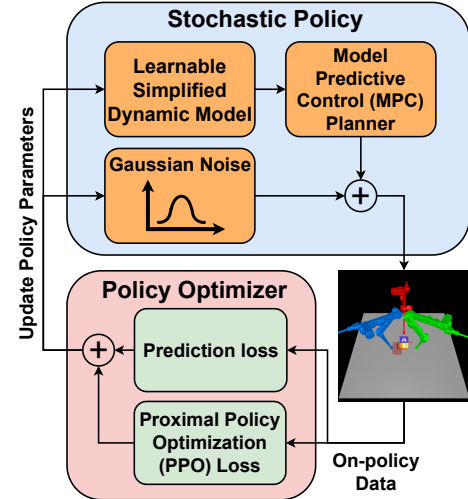


Fig. 1: The diagram demonstrates our proposed framework of learning simplified dynamic models for solving contact-rich manipulation tasks in low data regimes. Our framework proposes an iterative learning loop that consists of main components: a stochastic policy and a policy optimizer. **Top panel:** Using learned dynamic models under MPC scheme and Gaussian noise to construct the stochastic policy. **Bottom panel:** Combining PPO and prediction loss to optimize the policy parameters with the collected on-policy data.

to model-free RL, model-based RL remains data-intensive because conventional methods for learning deep dynamic models struggle to capture stiffness and multi-modal contact dynamics [15], [16]. Moreover, by adopting simpler model representations such as time-varying linear or Gaussian Process models, some works [17]–[19] demonstrate good performance with limited data.

The most recent work [20] shows great performance with only a few minutes of data by learning a linear complementarity system (LCS), a piecewise-affine and reduced-order representation of multi-contact dynamics, and combining it with an MPC planner. However, two main building blocks of this work, the dynamic model fitting and planning with the learned model, appear as two de-coupled optimization problems repeated over many cycles of learning. More specifically, ensuring better forward prediction capability of the learned dynamic model is necessary but might not be sufficient to achieve better task performance, leading to limitations of data efficiency and task performance. This issue is known as *objective mismatch* [21].

This paper presents LCS-RL, a low-dimensional RL framework, to address the above issue, aiming to further

\*Toyota Research Institute provided funds to support this work.

<sup>1</sup>The authors are with the GRASP Laboratory, University of Pennsylvania, Philadelphia, PA 19104, USA {xuanhien, posa}@seas.upenn.edu

enhance task performance and data efficiency. Particularly, the proposed framework directly bridges the dynamic model learning part to task performance optimization via a policy gradient algorithm.

### A. Contributions

- 1) We present LCS-RL, a novel framework that leverages the combination of RL and simple multi-contact models for solving contact-rich tasks. Specifically, our framework applies a reinforcement learning algorithm to directly maximize the task performance of simplified models in combination with the MPC planner.
- 2) We show that the proposed method consistently achieves higher task performance, up to 15% in three-fingered robot manipulation tasks with various objects compared to the prior methods [10], [20]. In addition, our method is data efficient as it can solve some dexterous manipulation tasks with 70% to 96% success rates using just under 30 minutes of data.
- 3) We also demonstrate that the learned LCS model of one object can be transferred to other objects, drastically improving data efficiency.

## II. RELATED WORK

### A. Differentiable MPC and Reinforcement Learning

Control policies formulated as differentiable MPC problems and optimized using backpropagation, either through imitation learning or RL loss, have been extensively studied [22]–[30]. Amos et al. [23], Xu et al. [30], and Jin et al. [26], [27] propose using either a differential set of Karush–Kuhn–Tucker (KKT) or Pontryagin’s maximum principle (PMP) conditions to ensure MPC problem differentiability, with limited experiments in low-dimensional tasks. Recent works [24], [25], [28] extend these methods to more complex robotics problems. Esfahani et al. [24] use a specialized Q-learning algorithm to learn MPC cost function parameters, effective in mobile robot tasks but critically requiring ground-truth dynamics. Wan et al. [28] focus on image-based tasks, introducing a differentiable sampling-based MPC policy to learn latent dynamics models, encoder, and Q-value predictors concurrently, but requiring substantial data. For data-efficient system identification and control in state-based contact-rich tasks, Saxena et al. [25] are the most relevant. They parameterize the dynamics model using a switching linear dynamical model with a contact-mode prediction function and construct a differentiable feedback controller (LQR), optimizing both its cost matrices and dynamics model parameters to match expert demonstrations.

### B. Main Baseline for Comparisons

Jin et al. [20] suggest that there exists a simple model that can adequately capture task-relevant contact dynamics, thereby enabling both high performance and real-time control for contact-rich manipulation. In particular, the authors propose to use a reduced-order hybrid model to represent and use the model predictive controller for planning. Since the model is far simpler than deep neural networks, much

less data is required for model learning. Their framework achieves high task performance with under 5 minutes of data. In this paper, we compare the task performance and data efficiency of our proposed method against this baseline in some dexterous manipulation tasks.

## III. BACKGROUNDS

### A. Linear Complementarity Systems

A discrete-time linear complementarity system (LCS) is a piecewise-affine system, where the state evolution is governed by linear dynamics in (1a) and a linear complementarity problem (LCP) in (1b).

$$\mathbf{x}_{t+1} = A\mathbf{x}_t + B\mathbf{u}_t + C\boldsymbol{\lambda}_t + \mathbf{d}, \quad (1a)$$

$$0 \leq \boldsymbol{\lambda}_t \perp D\mathbf{x}_t + E\mathbf{u}_t + F\boldsymbol{\lambda}_t + \mathbf{c} \geq 0. \quad (1b)$$

Here,  $\mathbf{x}_t \in \mathbb{R}^{n_x}$ ,  $\mathbf{u}_t \in \mathbb{R}^{n_u}$ , and  $\boldsymbol{\lambda}_t \in \mathbb{R}^{n_\lambda}$  are respectively the system state, action, and the complementarity variable at time step  $t$ . And,  $\mathbf{x}_{t+1} \in \mathbb{R}^{n_x}$  is the system state at the next time step  $t + 1$ . The symbol  $\perp$  denotes zero inner product or orthogonality of two vectors. Moreover, the matrix  $A \in \mathbb{R}^{n_x \times n_x}$  defines the autonomous dynamics and matrix  $B \in \mathbb{R}^{n_x \times n_u}$  captures the effect of actions on states. And, the matrix  $C \in \mathbb{R}^{n_x \times n_\lambda}$  and  $\mathbf{d} \in \mathbb{R}^{n_x}$  describe the effect of the contact forces and the constant forces acting on the state respectively. Other matrices  $D \in \mathbb{R}^{n_\lambda \times n_x}$ ,  $E \in \mathbb{R}^{n_\lambda \times n_u}$ ,  $F \in \mathbb{R}^{n_\lambda \times n_\lambda}$  and  $\mathbf{c} \in \mathbb{R}^{n_\lambda}$  altogether capture the relationship between states, actions, and contact forces. The LCS models are commonly used in modeling multi-contact robotics problems [31], [32].

### B. Learning Linear Complementarity Systems

Given a data buffer  $\mathcal{D}$  that contains some state transitions  $(\mathbf{x}_t, \mathbf{u}_t, \mathbf{x}_{t+1})$ , we can learn all matrix and vector parameters of an LCS model  $(A, B, C, \mathbf{d}, D, E, F, \mathbf{c})$  in (1) by using the gradient descent method with the violation-based loss, proposed by Jin et al. [33]

$$\begin{aligned} \mathcal{L}_{\text{vio}}^{\ominus} = & \min_{\boldsymbol{\lambda}_t \geq 0, \phi_t \geq 0} \frac{1}{2} \|A\mathbf{x}_t + B\mathbf{u}_t + C\boldsymbol{\lambda}_t + \mathbf{d} - \mathbf{x}_{t+1}\|^2 \\ & + \frac{1}{\xi} \left( \boldsymbol{\lambda}_t^\top \phi_t + \frac{1}{2\gamma} \|D\mathbf{x}_t + E\mathbf{u}_t + F\boldsymbol{\lambda}_t + \mathbf{c} - \phi_t\|^2 \right). \end{aligned} \quad (2)$$

Particularly, the loss  $\mathcal{L}_{\text{vio}}^{\ominus}$  itself is an optimization problem whose first and second terms specify the violation of the affine dynamics (1a) and the LCP constraint (1b), respectively. Under the condition  $0 < \gamma \leq \sigma_{\min}(F + F^T)$ , finding  $\boldsymbol{\lambda}_t$  to minimize the second term is equivalent to directly solving an LCP in (1b) for  $\boldsymbol{\lambda}_t$ , but poses a better-conditioned landscape for  $\mathcal{L}_{\text{vio}}^{\ominus}$ , thus enabling the identification of multi-modal and stiff dynamics. The hyper-parameter  $\xi > 0$  aims to balance two terms of the loss  $\mathcal{L}_{\text{vio}}^{\ominus}$ ; and  $\phi_t \in \mathbb{R}^{n_\lambda}$  is an introduced slack variable for the complementarity equation. Full explanations of the loss formulation and its hyper-parameters can be found in [33].

As proven in [33], using Envelope Theorem [34], we can analytically compute the gradient of the violation-based loss with respect to LCS parameters  $\frac{d\mathcal{L}_{\text{vio}}^{\ominus}}{d\Theta}$  without differentiating through the solution of the optimization problem.

---

**Algorithm 1:** Warm-up phase for optimizing LCS-MPC Policy

---

**Parameterization:** LCS-MPC policy parameters  $\theta$ ;  
**Hyper-parameters:** The number of warm-up iterations  $M$ ; the number of policy improvement steps  $N_p$ ; and learning rate  $\eta$ ;  
**Initialization:**  $\theta_0$ , empty data buffer  $\mathcal{D}$ ;

- 1 **for**  $k = 0, 1, \dots, M$  **do**
- 2     Collect  $N$  rollout trajectories by running the LCS-MPC stochastic policy  $\pi_{\theta_k}$  and add to  $\mathcal{D}$
- 3     **for**  $i \leftarrow 0$  **to**  $N_p$  **do**
- 4         Using data in  $\mathcal{D}$ , compute the gradient  $\frac{d\mathcal{L}_{\text{vio}}^{\theta_i}}{d\theta_i}$
- 5         Update  $\theta_{i+1} \leftarrow \theta_i - \eta \frac{d\mathcal{L}_{\text{vio}}^{\theta_i}}{d\theta_i}$
- 6     **end**
- 7 **end**
- 8 Save the final parameters  $\theta_M$  for the main phase

---

### C. Model Predictive Controller with LCS

Utilizing LCS to represent the dynamics model, one can construct a model predictive controller (MPC) as follows:

$$\begin{aligned}
 & \min_{\mathbf{u}_t, \mathbf{u}_{t+1}, \dots, \mathbf{u}_{t+H-1}} \sum_{k=t}^{t+H-1} \mathcal{C}(\mathbf{x}_k, \mathbf{u}_k) + \mathcal{C}_f(\mathbf{x}_{t+H}) \\
 & \text{s.t.} \quad \mathbf{x}_{k+1} = A\mathbf{x}_k + B\mathbf{u}_k + C\boldsymbol{\lambda}_k + \mathbf{d}, \\
 & \quad \mathbf{0} \leq \boldsymbol{\lambda}_k \perp D\mathbf{x}_k + E\mathbf{u}_k + F\boldsymbol{\lambda}_k + \mathbf{c} \geq \mathbf{0}, \\
 & \quad \mathbf{u}_{\min} \leq \mathbf{u}_k \leq \mathbf{u}_{\max}.
 \end{aligned} \tag{3}$$

where  $H$  is the planning horizon;  $\mathcal{C}$  and  $\mathcal{C}_f$  are the path and terminal cost functions. And,  $\mathbf{u}_{\min}$  and  $\mathbf{u}_{\max}$  are the lower and upper bounds of actions.

Given any initial state  $\mathbf{x}_t$ , we solve the LCS-MPC in (3) to plan a sequence of optimal actions  $[\mathbf{u}_t, \mathbf{u}_{t+1}, \dots, \mathbf{u}_{t+H-1}]$  that minimizes the total cost, then select the first action  $\mathbf{u}_t$  to apply on the robot and repeat the process in every time step in a receding horizon manner. To efficiently solve the LCS-MPC, we employ the direct trajectory optimization method [35], which simultaneously searches over trajectories of  $\mathbf{x}_{t:t+H}$ ,  $\mathbf{u}_{t:t+H-1}$ , and  $\boldsymbol{\lambda}_{t:t+H-1}$ , treating the LCS dynamics as a separate constraint for each time step. Also, we use the IPOPT solver [36] to solve such nonlinear problems.

### D. Proximal Policy Optimization

Proximal Policy Optimization (PPO) [5] is a policy gradient algorithm that focuses on determining how to make the most significant policy improvement using current data, all while avoiding excessive steps that could lead to performance collapse. Particularly, the PPO loss is defined as follows:

$$\mathcal{L}_{\text{PPO}}^{\theta} = -\frac{1}{|\mathcal{D}|T} \sum_{\tau \in \mathcal{D}} \sum_{t=0}^T \begin{cases} \max(h_t^{\theta}, 1 - \epsilon) \mathcal{A}_t & \text{if } \mathcal{A}_t < 0 \\ \min(h_t^{\theta}, 1 + \epsilon) \mathcal{A}_t & \text{if } \mathcal{A}_t \geq 0, \end{cases} \tag{4}$$

where  $\mathcal{D}$  and  $|\mathcal{D}|$  are the data buffer and its size, that buffer consists of on-policy rollout trajectories  $\tau$ , and  $T$  is the length of trajectories. There are two key quantities in (4): the ratio  $h_t^{\theta} = \frac{\pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t)}{\pi_{\theta_{\text{old}}}(\mathbf{u}_t | \mathbf{x}_t)}$  and the truncated version of

---

**Algorithm 2:** Main phase for optimizing LCS-MPC Policy using the PPO algorithm

---

**Parameterization:** LCS-MPC policy parameters  $\theta$  and value function  $\mathcal{V}_{\phi}$ ;  
**Hyper-parameters:** Total number of iterations  $K$ , the number of policy improvement steps  $\bar{N}_p$ ; Learning rate  $\bar{\eta}$  for the policy optimization; Loss weighting parameter  $\beta$  in (7); Discount factor  $\gamma$  and parameter  $\zeta$  for computing the advantage values;  
**Initialization:**  $\theta_M$  obtained from the warm-up phase,  $\phi_0$ , and empty data buffer  $\mathcal{D}$ ;

- 1 **for**  $k = 0, 1, \dots, K$  **do**
- 2     Empty buffer  $\mathcal{D}$ , collect  $\bar{N}$  new trajectories by running the LCS-MPC policy  $\pi_{\theta_k}$ , and add to  $\mathcal{D}$
- 3     For each trajectory, compute the generalized advantage value  $\mathcal{A}_t$  as in (5), then the bootstrapped total reward  $R_t = \mathcal{A}_t + \mathcal{V}_{\phi_k}(\mathbf{x}_t)$
- 4     **for**  $i \leftarrow 0$  **to**  $\bar{N}_p$  **do**
- 5         Compute the combined loss gradient  $\frac{d\mathcal{L}_c^{\theta_i}}{d\theta_i}$
- 6         Update  $\theta_{i+1} \leftarrow \theta_i - \bar{\eta} \frac{d\mathcal{L}_c^{\theta_i}}{d\theta_i}$
- 7     **end**
- 8     Fit value function  $\mathcal{V}_{\phi}$  by performing regression with mean-square error
- 9     
$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}|T} \sum_{\tau \in \mathcal{D}} \sum_{t=0}^T (\mathcal{V}_{\phi}(\mathbf{x}_t) - R_t)^2$$

---

generalized advantage function  $\mathcal{A}_t$  [37]. Here, the ratio  $h_t^{\theta}$  indicates how much the new policy differs from the old one. The scalar  $\epsilon$  defines the bounds of  $h_t^{\theta}$ , which are often referred to as the trust region of policy improvements. In addition,  $\mathcal{A}_t$  guides the policy search by measuring whether a certain action is a good or bad decision within a given state. The detailed expression of  $\mathcal{A}_t$  is given below

$$\begin{aligned}
 \mathcal{A}_t &= \delta_t + (\gamma\zeta)\delta_{t+1} + (\gamma\zeta)^2\delta_{t+2} \dots + (\gamma\zeta)^{T-t+1}\delta_{T-1}, \\
 & \text{with } \delta_t = r_t + \gamma\mathcal{V}_{\phi}(\mathbf{x}_{t+1}) - \mathcal{V}_{\phi}(\mathbf{x}_t),
 \end{aligned} \tag{5}$$

where  $\gamma$  is the discount factor and  $\zeta$  is a hyper-parameter that controls the bias-variance tradeoff of the estimation. The reward  $r_t$  is obtained by executing action  $\mathbf{u}_t$  at state  $\mathbf{x}_t$ . The learned value function  $\mathcal{V}_{\phi}(\mathbf{x}_t)$  and  $\mathcal{V}_{\phi}(\mathbf{x}_{t+1})$  estimate the expected total rewards if we follow the current policy from state  $\mathbf{x}_t$  and  $\mathbf{x}_{t+1}$  till the end of trajectories.

## IV. PRACTICAL ALGORITHM

In this section, we introduce our framework, LCS-RL, that utilizes a reinforcement learning algorithm, here we use PPO, to optimize an LCS dynamic model (in combination with model predictive control) for solving contact-rich tasks.

First, we formulate a stochastic policy, called the LCS-MPC stochastic policy, by adding Gaussian noise to the output of the LCS-MPC planner in (3). In other words, the LCS-MPC policy is directly parameterized by the LCS model. Then, we use the combination of the PPO loss and

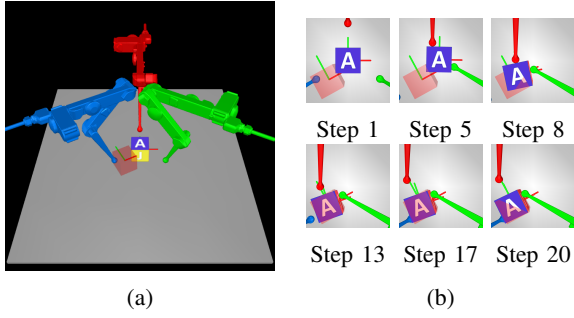


Fig. 2: TriFinger dexterous manipulation tasks. (a) shows the simulation environment that is constructed using MuJoCo physics engine [38]. In this task, the three fingers need to push the cube towards a random target pose, visualized by the red transparent cube. (b) is an example of a rollout trajectory that demonstrates how the fingers approach, make, and break contacts to reposition the cube.

the violation-based loss given in (7) to improve the task performance of that policy.

We address the poor data efficiency of the PPO algorithm by leveraging data-efficient model learning at the start and then transitioning to PPO when in a good neighborhood. Therefore, our framework consists of two phases: the warm-up phase and the main phase. In the warm-up phase, we follow the algorithm proposed in [20], solely employing the violation-based loss (2) to quickly learn the parameters of the LCS model that can achieve good task performance. Subsequently, we use the learned LCS model to accelerate the main phase, where we start using the PPO algorithm. In practice, we find that having the warm-up phase leads to more stable and progressive training than involving PPO right from the beginning. While theoretically, we could merge the two phases and switch the loss upon transition, it is simpler to keep them separated due to different hyper-parameters required for each phase. Details of the warm-up phase and main phase are provided in Algorithm 1 and 2.

#### A. LCS-MPC Stochastic Policy

The LCS-MPC policy is the probability density of an action distribution associated with the current state  $\mathbf{x}_t$ :

$$\pi_{\theta}(\mathbf{u}_t|\mathbf{x}_t) = \frac{\exp(-\frac{1}{2}(\mathbf{u}_t - \mu_{\Theta}(\mathbf{x}_t))^T \Sigma^{-1}(\mathbf{u}_t - \mu_{\Theta}(\mathbf{x}_t)))}{(2\pi)^{n_u/2} \det(\Sigma)^{1/2}}, \quad (6)$$

where  $\mu_{\Theta} \in \mathbb{R}^{n_u}$  is the optimal output of the LCS-MPC planner, and  $\Sigma \in \mathbb{R}^{n_u \times n_u}$  is the covariance matrix that indicates the noise magnitude. Here,  $\Theta = (A, B, C, \mathbf{d}, D, E, F, \mathbf{c})$  is actually the LCS parameters. And,  $\theta = [\Theta, \Sigma]$  denotes the joint vector of the policy's learnable parameters. The added noise encourages exploration, helping to avoid low-quality local minima. Typically, the noise magnitude is large initially, gradually decreasing as the policy exploits acquired knowledge for better task performance.

#### B. Loss for Optimizing LCS Model

In our framework, we employ two types of loss: violation-based and PPO. The violation-based loss improves the forward prediction capability of the LCS model, while the PPO

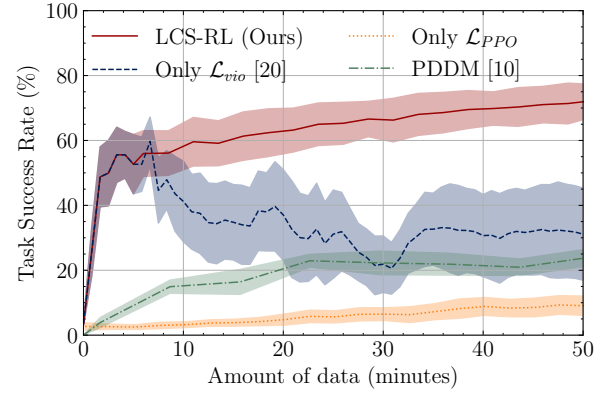


Fig. 3: Learning curves of the TriFinger Moving Cube task. The red, blue, orange, and green lines show the average task success rate of our proposed method, the prior method [20], a method that uses PPO without a warm-up phase, and PDDM [10] respectively. At the beginning of the training, our method and the prior method [20] share the same performance since the same algorithm is used. However, the transition occurs after collecting 6 minutes of data, when our method switches to fully employ the PPO algorithm. Shaded regions indicate normal t-score 95% confidence intervals.

loss directly enhances the task performance of the LCS-MPC planner. Introducing a hyper-parameter  $\beta \in [0, 1]$  allows us to balance the contributions of these losses, resulting in the combined loss:

$$\mathcal{L}_c^{\theta} = \beta \mathcal{L}_{\text{PPO}}^{\theta} + (1 - \beta) \mathcal{L}_{\text{vio}}^{\theta}, \quad (7)$$

In order to optimize the parameters of LCS-MPC stochastic policy, one must compute the gradient of the PPO loss and the violation-based loss with respect to the policy parameters. In section III-B we have already mentioned the method for calculating  $\frac{d\mathcal{L}_{\text{vio}}^{\theta}}{d\theta}$ . Moving forward, we will illustrate the process for computing  $\frac{d\mathcal{L}_{\text{PPO}}^{\theta}}{d\theta}$ .

#### C. Gradient of PPO Loss

The original PPO algorithm parameterizes policies via deep neural networks [5], thus computing the gradient of the PPO loss over policy parameters can be simply done via automatic differentiation. However, it does not apply to our case since our policy is actually an MPC planner. Differentiating through the MPC requires special treatment.

Given the explicit form of PPO loss in (4), using chain rule, the gradient with respect to  $\theta$  can be computed

$$\frac{d\mathcal{L}_{\text{PPO}}^{\theta}}{d\theta} = \frac{-1}{|\mathcal{D}|T} \sum_{\tau \in \mathcal{D}} \sum_{t=0}^T \begin{cases} \frac{dh_t^{\theta}}{d\pi_{\theta}} \frac{d\pi_{\theta}}{d\theta} \mathcal{A}_t & \text{if } h_t^{\theta} \geq 1 - \epsilon; \mathcal{A}_t < 0 \\ 0 & \text{if } h_t^{\theta} < 1 - \epsilon; \mathcal{A}_t < 0 \\ \frac{dh_t^{\theta}}{d\pi_{\theta}} \frac{d\pi_{\theta}}{d\theta} \mathcal{A}_t & \text{if } h_t^{\theta} \leq 1 + \epsilon; \mathcal{A}_t \geq 0 \\ 0 & \text{if } h_t^{\theta} > 1 + \epsilon; \mathcal{A}_t \geq 0. \end{cases} \quad (8)$$

Due to the clipping effect of the PPO loss, the gradients are zero when the improvement steps of the PPO policy  $h_t^{\theta}$  are outside of trust region  $[1 - \epsilon, 1 + \epsilon]$ . Hence, we are left to compute the gradients if  $h_t^{\theta}$  stays within the trust region. To compute (8), one must compute  $\frac{dh_t^{\theta}}{d\pi_{\theta}}$  and  $\frac{d\pi_{\theta}}{d\theta}$ . While it

Object	Peak Success Rate (%)		Final Success Rate (%)	
	Only $\mathcal{L}_{\text{vio}}^{\theta}$	LCS-RL (Ours)	Only $\mathcal{L}_{\text{vio}}^{\theta}$	LCS-RL (Ours)
Sugar Box	87.5 ± 9.2	<b>95.9 ± 2.4</b>	44.3 ± 26.2	<b>95.9 ± 2.4</b>
Fish Can	59.1 ± 7.3	<b>69.9 ± 8.1</b>	38.8 ± 13.1	<b>69.9 ± 8.1</b>
Mug	32.5 ± 8.4	<b>44.6 ± 3.9</b>	10.0 ± 6.3	<b>44.6 ± 3.9</b>
Wrench	45.5 ± 8.7	<b>60.0 ± 6.5</b>	11.5 ± 8.6	<b>59.5 ± 6.6</b>
Clamp	24.7 ± 5.3	<b>39.3 ± 9.7</b>	6.1 ± 5.6	<b>38.7 ± 10.2</b>
Banana	28.4 ± 5.8	<b>35.5 ± 5.3</b>	7.5 ± 6.6	<b>35.5 ± 5.3</b>

TABLE I: Comparison task success rates between our method and prior method [20] on diverse objects.

is straightforward to evaluate  $\frac{dh_t}{d\pi_{\theta}}$ , computing  $\frac{d\pi_{\theta}}{d\theta}$  requires differentiation of the optimal actions of the MPC  $\mu_{\Theta}(x_t)$  with respect to its parameters  $\Theta$ . We compute this derivative via perturbations of KKT conditions [39] with details given in [40]. Also, note that the MPC problem is not always classically differentiable (e.g. when strict complementarity does not hold in the KKT conditions), but we have not found this problematic in practice.

## V. EXPERIMENTS AND RESULTS

In this section, we will verify our proposed framework on the three-fingered robotic hand manipulation task that was first proposed by [20] (see Fig. 2). We call it the TriFinger Moving Cube task. In the first experiment, we show a comparison of the task performance of the LCS model trained by our method and prior methods. Next, we replace the cube with other objects that have more complex shapes and repeat the same experiment. Lastly, we demonstrate that the data efficiency of our framework can be greatly improved via transfer learning. In order to guarantee statistically meaningful results, for each experiment, we have 10 runs with 10 random seeds. Also, we compute the task success rate by evaluating the learned models with 1000 random goal poses and aggregate results. All videos and codes are available at <https://sites.google.com/view/lcs-rl>.

### A. TriFinger Moving Cube Task

In this task, a TriFinger robot aims to align a 6 cm-sized cube with random goal poses on a planar surface. Each episode comprises up to 20 steps, each lasting 0.1 seconds.

1) *States and Actions*: We define the system state

$$\mathbf{x} = [\mathbf{p}_{\text{cube}}, \alpha_{\text{cube}}, \mathbf{p}_{\text{fingertips}}] \in \mathbb{R}^9, \quad (9)$$

where  $\mathbf{p}_{\text{cube}} \in \mathbb{R}^2$  is the xy position of the cube;  $\alpha_{\text{cube}} \in \mathbb{R}$  is the rotation angle around the z (vertical) axis; and  $\mathbf{p}_{\text{fingertips}} \in \mathbb{R}^6$  are the xy positions of three fingertips. We define the actions as deviations from the current positions of three fingertips in the Cartesian space. In addition, we impose safety limits on actions (element-wise) to constrain how far fingertips can move in one time step

$$\begin{aligned} \mathbf{u} &= \Delta \mathbf{p}_{\text{fingertips}} \in \mathbb{R}^6, \\ u_i &\in [-0.015, 0.015] \text{ m}. \end{aligned} \quad (10)$$

We employ operational space control (OSC) [41] in the lower-level controller to map action  $\mathbf{u}$  to the joint torque of each finger. We also utilize the OSC controller to maintain fingertips at a constant height as this TriFinger task involves only planar manipulation.

2) *Task Space*: We use the same bounds for task space as in [20], from which the goal poses are uniformly sampled

$$\begin{aligned} [-0.06, -0.06]^T \text{ m} &\leq \mathbf{p}_{\text{goal}} \leq [0.06, 0.06]^T \text{ m}, \\ -0.5 \text{ rads} &\leq \alpha_{\text{goal}} \leq 0.5 \text{ rads}. \end{aligned} \quad (11)$$

3) *Task Success Criteria*: When the cube pose is near the goal pose and within some tolerances, we can consider that the task is successfully completed. The goal tolerance values are selected to establish the right level of difficulty as too stringent tolerances make the task impossible to solve. We follow some previous works on TriFinger tasks [42], [43] to set the tolerances as follows:

$$\begin{aligned} \|\mathbf{p}_{\text{cube}}^* - \mathbf{p}_{\text{goal}}\| &\leq 0.02 \text{ m}, \\ \|\alpha_{\text{cube}}^* - \alpha_{\text{goal}}\| &\leq 0.2 \text{ rads}, \end{aligned} \quad (12)$$

where  $\mathbf{p}_{\text{cube}}^*$  and  $\alpha_{\text{cube}}^*$  are the xy position and orientation of the cube at the last time step  $T$ ;  $\mathbf{p}_{\text{goal}} \in \mathbb{R}^2$  and  $\alpha_{\text{goal}} \in \mathbb{R}$  together specify the goal pose.

4) *Cost function for the LCS-MPC*: We utilize the same cost function for the LCS-MPC as in the prior work [20]:

$$\begin{aligned} \mathcal{J} &= \sum_{k=t}^{t+H-1} \mathcal{C}(\mathbf{x}_k, \mathbf{u}_k) + \mathcal{C}_f(\mathbf{x}_{t+H}), \\ \mathcal{C} &= 10.0 \|\mathbf{p}_{\text{fingertips}} - \mathbf{p}_{\text{cube}}\|^2 + 200.0 \|\mathbf{p}_{\text{cube}} - \mathbf{p}_{\text{goal}}\|^2 \\ &\quad + 0.3 (\alpha_{\text{cube}} - \alpha_{\text{goal}})^2 + 200.0 \|\mathbf{u}\|^2, \\ \mathcal{C}_f &= 6.0 \|\mathbf{p}_{\text{fingertips}} - \mathbf{p}_{\text{cube}}\|^2 + 200.0 \|\mathbf{p}_{\text{cube}} - \mathbf{p}_{\text{goal}}\|^2 \\ &\quad + 1.5 (\alpha_{\text{cube}} - \alpha_{\text{goal}})^2, \end{aligned} \quad (13)$$

5) *Reward function for PPO*: We use both dense and sparse reward functions for the PPO algorithm. The dense reward function  $r_t(\mathbf{x}_t, \mathbf{u}_t)$  is simply the negation of the cost function  $\mathcal{C}$  in (13). This choice of reward function ensures that both PPO and the MPC planner of the stochastic policy align in the same direction toward task completion.

At the end of the rollout trajectory, we add a negative sparse reward to penalize for not completing the task:

$$r(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}) = -10.0 \times (1 - \text{is\_task\_completed}), \quad (14)$$

In practice, we find that sparse reward helps to accelerate the training significantly.

### B. Results of the TriFinger Moving Cube Task

To demonstrate the effectiveness of LCS-RL, we compare against the prior method [20], which trains purely on  $\mathcal{L}_{\text{vio}}^{\theta}$ , against PPO without a warm-up phase. We also compare against a state-of-the-art model-based RL approach PDDM [10] to demonstrate the utility of simple, non-smooth models over deep neural networks for manipulation. Note that in the main phase of our method, we set  $\beta = 1.0$  for the combined loss in (7), meaning that only the PPO loss is used. The results are shown in Fig.3.

When utilizing only the violation-based loss, the mean success rate peaks at 55% after 7 minutes of data, then fluctuates and decreases as more data is collected, and finally stops at approximately 30%. In contrast, starting with the same performance at 6 minutes of data, our method improves

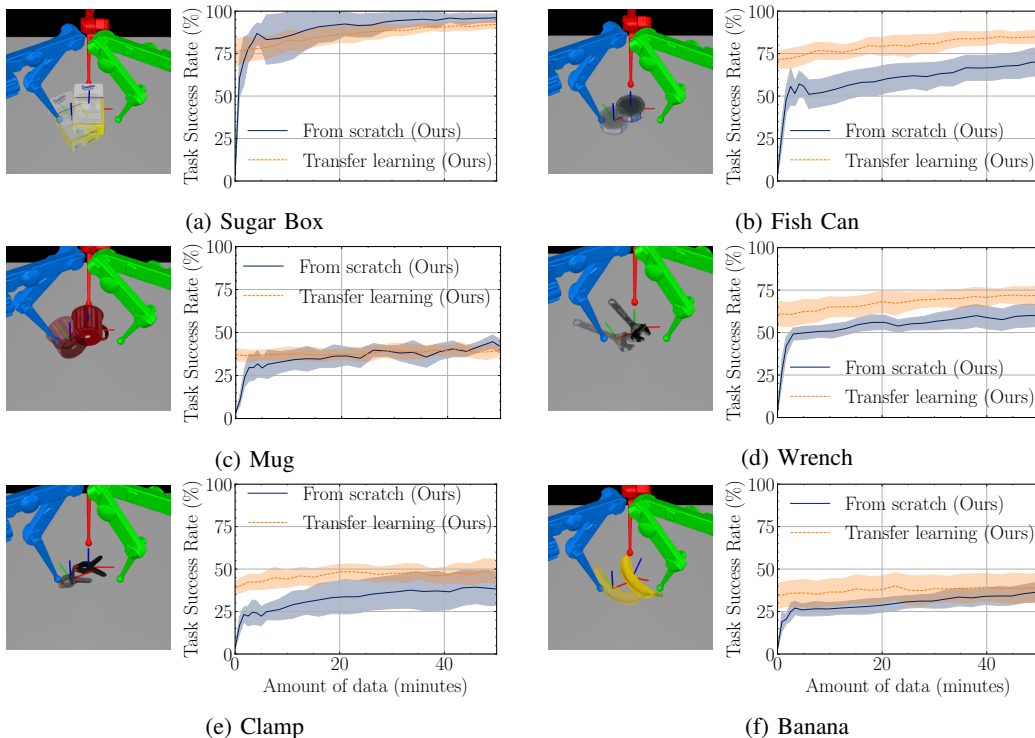


Fig. 4: Comparison task performance between learning the LCS model from scratch and pre-trained models (obtained from training with the TriFinger Moving Cube task) on the YCB objects using our LCS-RL framework.

the task performance throughout the training, reaching 65% of success rate after 25 minutes of data and 71.4% at the end of the training. Since the LCS models have limited expressiveness power, even if we optimize LCS models for better capability of forward prediction, this capability might not be optimally assigned to regions of state space where accuracy is needed for task performance. As a result, the overall task performance might drop significantly. Our method does not suffer from that issue because the only objective of PPO is encouraging the policy, to repeat good trajectories and avoid bad trajectories.

In addition, the PPO-only method and PDDM [10] have the lowest task performances throughout the training, achieving merely 10% and 20% for the final task success rate.

### C. TriFinger Manipulating Diverse Objects

We run a set of experiments on the TriFinger Moving Object task, which is similar to the TriFinger Moving Cube task, but the cube is replaced by other objects with non-convex, highly intricate shapes. Those objects, including sugar box, fish can, mug, wrench, clamp, and banana, are selected from the YCB object and model set [44]. As seen from Table I, our method consistently outperforms the prior method in [20] given the same amount of data, gaining from 8% (sugar box) to 15% (clamp) higher task success rate.

### D. Transfer Learning

To illustrate the transfer learning capabilities of our LCS-RL framework, we employ the LCS model initially trained on the TriFinger Moving Cube task as the starting point for training on other objects. The results in Fig. 4 show

that our LCS-RL framework is highly suitable for transfer learning. Particularly, we can observe that transfer learning significantly accelerates the training, yielding even higher final task success rates in all objects (except for the sugar box), compared to the training from scratch model.

## VI. CONCLUSIONS

In conclusion, we present LCS-RL, a novel approach that leverages a reinforcement learning algorithm to directly maximize the task performance of the LCS model in combination with the MPC planner. We demonstrate that the proposed method attains higher task performance and greater sample efficiency compared to prior methods in TriFinger robot tasks involving pushing and rotating various objects. In addition, we show that our method is highly suitable for transfer learning, which further helps to improve data efficiency.

Our framework is not limited to only the PPO algorithm since any RL algorithms can be incorporated. Thus, one direction for future work is to explore other RL algorithms and employ them in our framework. There are off-policy RL algorithms such as TD3 [6] or SAC [7], known for better data efficiency when compared to on-policy algorithms [45]. Nevertheless, this advantage may not be as evident in situations with limited data.

Lastly, we observe the limitation of using the LCS model for representing system dynamic models, especially with complex geometries like bananas or clamps. We aim to investigate alternative structured models simpler than neural networks yet exhibiting nonlinear components, although with a trade-off in data efficiency. One potential candidate is the Nonlinear Complementarity System (NCS) model.

## REFERENCES

- [1] T. Marcucci and R. Tedrake, "Warm Start of Mixed-Integer Programs for Model Predictive Control of Hybrid Systems," *IEEE Transactions on Automatic Control*, vol. 66, pp. 2433–2448, 2019.
- [2] D. Frick, A. Georghiou, J. L. Jerez, A. Domahidi, and M. Morari, "Low-complexity method for hybrid MPC with local guarantees," *SIAM Journal on Control and Optimization*, vol. 57, no. 4, pp. 2328–2361, 2019.
- [3] A. Aydinoglu, A. Wei, and M. Posa, "Consensus Complementarity Control for Multi-Contact MPC," *arXiv preprint arXiv:2304.11259*, 2023.
- [4] A. Aydinoglu and M. Posa, "Real-Time Multi-Contact Model Predictive Control via ADMM," in *2022 International Conference on Robotics and Automation (ICRA)*. Philadelphia, PA, USA: IEEE, 2022, pp. 3414–3421.
- [5] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [6] S. Fujimoto, H. Hoof, and D. Meger, "Addressing Function Approximation Error in Actor-Critic Methods," in *Proceedings of the 35th International Conference on Machine Learning*. PMLR, 2018, pp. 1587–1596.
- [7] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor," in *Proceedings of the 35th International Conference on Machine Learning*, vol. 80. PMLR, 2018, pp. 1861–1870.
- [8] A. Nagabandi, G. Kahn, R. S. Fearing, and S. Levine, "Neural Network Dynamics for Model-Based Deep Reinforcement Learning with Model-Free Fine-Tuning," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 7559–7566.
- [9] K. Chua, R. Calandra, R. McAllister, and S. Levine, "Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models," in *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [10] A. Nagabandi, K. Konoglie, S. Levine, and V. Kumar, "Deep Dynamics Models for Learning Dexterous Manipulation," in *Conference on Robot Learning*. PMLR, 2020, pp. 1101–1112.
- [11] A. S. Morgan, D. Nandha, G. Chalvatzaki, C. D'Eramo, A. M. Dollar, and J. Peters, "Model Predictive Actor-Critic: Accelerating Robot Skill Acquisition with Deep Reinforcement Learning," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 6672–6678.
- [12] F. Ebert, C. Finn, S. Dasari, A. Xie, A. Lee, and S. Levine, "Visual Foresight: Model-Based Deep Reinforcement Learning for Vision-Based Robotic Control," *arXiv:1812.00568 [cs]*, 2018.
- [13] M. Zhang, S. Vikram, L. Smith, P. Abbeel, M. J. Johnson, and S. Levine, "SOLAR: Deep Structured Representations for Model-Based Reinforcement Learning," in *Proceedings of the 36th International Conference on Machine Learning*, vol. 97. PMLR, 2019, pp. 7444–7453.
- [14] R. Ghugare, H. Bharadhwaj, B. Eysenbach, S. Levine, and R. Salakhutdinov, "Simplifying model-based RL: Learning representations, latent-space models, and policies with one objective," in *The Eleventh International Conference on Learning Representations*, 2023.
- [15] M. Parmar, M. Halm, and M. Posa, "Fundamental Challenges in Deep Learning for Stiff Contact Dynamics," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Prague, Czech Republic: IEEE, 2021, pp. 5181–5188.
- [16] B. Bianchini, M. Halm, N. Matni, and M. Posa, "Generalization Bounded Implicit Learning of Nearly Discontinuous Functions," in *Proceedings of The 4th Annual Learning for Dynamics and Control Conference (L4DC)*, ser. Proceedings of Machine Learning Research, vol. 168, 2022, pp. 1112–1124.
- [17] V. Kumar, E. Todorov, and S. Levine, "Optimal control with learned local models: Application to dexterous manipulation," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. Stockholm, Sweden: IEEE, 2016, pp. 378–383.
- [18] S. Levine and V. Koltun, "Guided Policy Search," in *Proceedings of the 30th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, S. Dasgupta and D. McAllester, Eds., vol. 28. Atlanta, Georgia, USA: PMLR, 2013, pp. 1–9.
- [19] M. P. Deisenroth and C. E. Rasmussen, "PILCO: A Model-Based and Data-Efficient Approach to Policy Search," in *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ser. ICML'11. Madison, WI, USA: Omnipress, 2011, pp. 465–472.
- [20] W. Jin and M. Posa, "Task-driven hybrid model reduction for dexterous manipulation," *IEEE Transactions on Robotics (TRO)*, vol. 40, pp. 1774–1794, Jan. 2024.
- [21] N. Lambert, B. Amos, O. Yadan, and R. Calandra, "Objective Mismatch in Model-based Reinforcement Learning," *arXiv preprint arXiv:2002.04523*, 2021.
- [22] M. Okada, L. Rigazio, and T. Aoshima, "Path Integral Networks: End-to-End Differentiable Optimal Control," *arXiv preprint arXiv:1706.09597*, 2017.
- [23] B. Amos, I. D. J. Rodriguez, J. Sacks, B. Boots, and J. Z. Kolter, "Differentiable MPC for End-to-end Planning and Control," in *Advances in Neural Information Processing Systems*, 2019.
- [24] H. N. Esfahani, A. B. Kordabad, and S. Gros, "Approximate Robust NMPC using Reinforcement Learning," in *2021 European Control Conference (ECC)*, Rotterdam, Netherlands, 2021.
- [25] S. Saxena, A. LaGrassa, and O. Kroemer, "Learning reactive and predictive differentiable controllers for switching linear dynamical models," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 7563–7569.
- [26] W. Jin, Z. Wang, Z. Yang, and S. Mou, "Pontryagin differentiable programming: An end-to-end learning and control framework," in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 7979–7992.
- [27] W. Jin, S. Mou, and G. J. Pappas, "Safe pontryagin differentiable programming," in *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc., 2021, pp. 16034–16050.
- [28] W. Wan, Y. Wang, Z. Erickson, and D. Held, "DiffTop: Differentiable trajectory optimization for deep reinforcement and imitation learning," *arXiv preprint arXiv:2402.05421*, 2024.
- [29] L. Pineda, T. Fan, M. Monge, S. Venkataraman, P. Sodhi, R. T. Chen, J. Ortiz, D. DeTone, A. Wang, S. Anderson, J. Dong, B. Amos, and M. Mukadam, "Theseus: A Library for Differentiable Nonlinear Optimization," *Advances in Neural Information Processing Systems*, 2022.
- [30] M. Xu, T. L. Molloy, and S. Gould, "Revisiting implicit differentiation for learning problems in optimal control," in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [31] D. Stewart and J. Trinkle, "An implicit time-stepping scheme for rigid body dynamics with coulomb friction," in *Proceedings 2000 ICRA. Millennium Conference. IEEE International Conference on Robotics and Automation. Symposia Proceedings (Cat. No.00CH37065)*, vol. 1, 2000, pp. 162–169 vol.1.
- [32] A. Aydinoglu, P. Sieg, V. M. Preciado, and M. Posa, "Stabilization of Complementarity Systems via Contact-Aware Controllers," *IEEE Transactions on Robotics*, vol. 38, no. 3, pp. 1735–1754.
- [33] W. Jin, A. Aydinoglu, M. Halm, and M. Posa, "Learning Linear Complementarity Systems," in *Proceedings of The 4th Annual Learning for Dynamics and Control Conference (L4DC)*. PMLR, 2022, p. 13.
- [34] S. N. Afriat, "Theory of maxima and the method of lagrange," *SIAM Journal on Applied Mathematics*, vol. 20, no. 3, pp. 343–357, 1971.
- [35] M. Posa, C. Cantu, and R. Tedrake, "A direct method for trajectory optimization of rigid bodies through contact," *The International Journal of Robotics Research*, vol. 33, no. 1, pp. 69–81, 2014.
- [36] A. Wächter and L. T. Biegler, "On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming," *Mathematical Programming*, vol. 106, no. 1, pp. 25–57, 2006.
- [37] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-Dimensional Continuous Control Using Generalized Advantage Estimation," in *4th International Conference on Learning Representations (ICLR)*, 2016.
- [38] E. Todorov, T. Erez, and Y. Tassa, "MuJoCo: A physics engine for model-based control," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 5026–5033.
- [39] H. W. Kuhn and A. W. Tucker, "Nonlinear programming," in *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability, 1950*. Berkeley and Los Angeles: University of California Press, 1951, pp. 481–492.
- [40] C. Büskens and H. Maurer, "Sensitivity analysis and real-time control of nonlinear optimal control systems via nonlinear programming methods," in *Variational Calculus, Optimal Control and Applications: International Conference in Honour of L. Bittner and R. Klötzler*, Basel: Birkhäuser Basel, 1998.

- [41] O. Khatib, "A unified approach for motion and force control of robot manipulators: The operational space formulation." *IEEE J. Robotics Autom.*, vol. 3, pp. 43–53, 1987.
- [42] A. Allshire, M. Mittal, V. Lodaya, V. Makoviychuk, D. Makoviichuk, F. Widmaier, M. Wüthrich, S. Bauer, A. Handa, and A. Garg, "Transferring dexterous manipulation from gpu simulation to a remote real-world trifinger," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 11 802–11 809.
- [43] N. Funk, C. Schaff, R. Madan, T. Yoneda, J. U. De Jesus, J. Watson, E. K. Gordon, F. Widmaier, S. Bauer, S. S. Srinivasa, T. Bhattacharjee, M. R. Walter, and J. Peters, "Benchmarking Structured Policies and Policy Optimization for Real-World Dexterous Object Manipulation," *IEEE Robotics and Automation Letters*, vol. 7, no. 1, pp. 478–485, 2022.
- [44] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollar, "The YCB object and Model set: Towards common benchmarks for manipulation research," in *2015 International Conference on Advanced Robotics (ICAR)*. Istanbul, Turkey: IEEE, 2015, pp. 510–517.
- [45] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2019.