

# Towards Human-Level Learning of Complex Physical Puzzles

Kei Ota<sup>1</sup>, Devesh K. Jha<sup>2</sup>, Diego Romeres<sup>2</sup>, Jeroen van Baar<sup>2</sup>, Kevin A. Smith<sup>3</sup>, Takayuki Semitsu<sup>1</sup>,  
Tomoaki Oiki<sup>1</sup>, Alan Sullivan<sup>2</sup>, Daniel Nikovski<sup>2</sup>, and Joshua B. Tenenbaum<sup>3</sup>

<sup>1</sup>Mitsubishi Electric, <sup>2</sup>Mitsubishi Electric Research Labs, <sup>3</sup>Massachusetts Institute of Technology

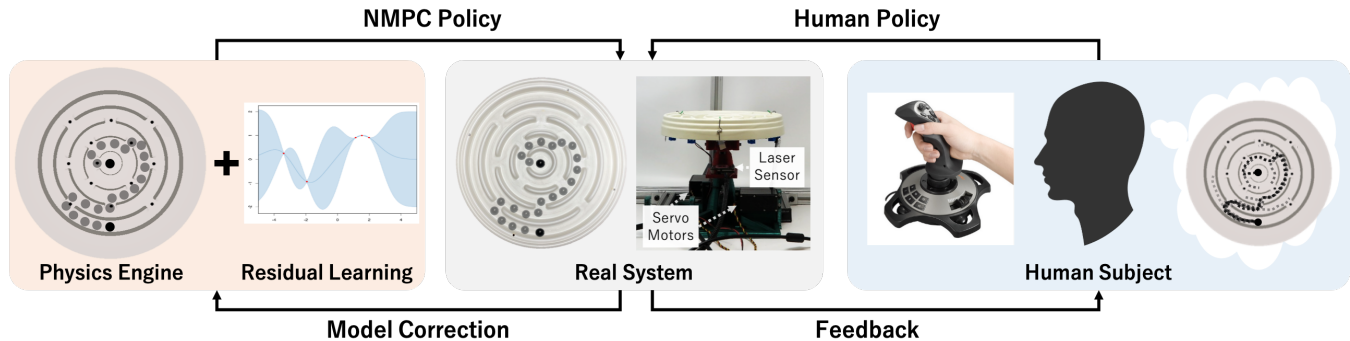


Fig. 1: Humans can learn to solve complex physical tasks with very little interaction with a system. Studies in cognitive science suggests that people have internal models of physics which are calibrated as they interact with new systems. In this paper, we train a reinforcement learning agent that initializes a policy with a general purpose physics engine, then we correct its model of dynamics using parameter estimation and residual model learning. We compare this learning method on a marble-in-a-maze puzzle and compare its behavior with how people perform in this environment.

**Abstract**—Humans quickly solve tasks in novel systems with complex dynamics, without requiring much interaction. While deep reinforcement learning algorithms have achieved tremendous success in many complex tasks, these algorithms need a large number of samples to learn meaningful policies. In this paper, we present a task for navigating a marble to the center of a circular maze. While this system is very intuitive and easy for humans to solve, it can be very difficult and inefficient for standard reinforcement learning algorithms to learn meaningful policies. We present a model that learns to move a marble in the complex environment within minutes of interacting with the real system. Learning consists of initializing a physics engine with parameters estimated using data from the real system. The error in the physics engine is then corrected using Gaussian process regression, which is used to model the residual between real observations and physics engine simulations. The physics engine equipped with the residual model is then used to control the marble in the maze environment using a model-predictive feedback over a receding horizon. We contrast the learning behavior against the time taken by humans to solve the problem to show comparable behavior. To the best of our knowledge, this is the first time that a hybrid model consisting of a full physics engine along with a statistical function approximator has been used to control a complex physical system in real-time using nonlinear model-predictive control (NMPC). Codes for the simulation environment can be downloaded here<sup>1</sup>. A video describing our method could be found here<sup>2</sup>.

## I. INTRODUCTION

People have remarkable capabilities to interact with the world in flexible and generalizable ways. With very little

effort, they can figure out how to use novel objects to accomplish their goals, or manipulate existing objects in new ways [1]. Artificial intelligence has long had the goal of designing robotic agents that can interact with the physical world in these human-like ways [2], [3]. Some of this work uses model-based control methods that form plans based on predefined models of the world dynamics. However, these systems require the accurate dynamics models, but even the best simulators diverge from the real world in some ways. Other recent work in machine learning has treated this as a reinforcement learning problem, assuming that their agents will learn a model of the world dynamics in tandem with control policies [4], [5], [6]. However, whereas these systems perform well in the scenarios they were trained to solve, they often fail to learn a model that is as flexible or efficient as people have [7], [8].

Our aim in this paper is to combine the best of both methodologies: our system uses nonlinear model predictive control with a predefined model of dynamics at its core, but updates that model by learning residuals between predictions and real-world observations via physical parameter estimation and Gaussian process regression [9]. This approach is inspired by work from cognitive science that suggests people have internal models of physics that are well calibrated to the world [10], [11], and that they use these models to learn how to use new objects to accomplish novel goals in just a handful of interactions [12]. In this way, we hope to attain human levels of physical control, while also achieving human levels of sample efficiency. A broad idea of the proposed approach

<sup>1</sup><https://www.merl.com/research/license/CME>

<sup>2</sup><https://youtu.be/xaxNCXBovpc>

is provided in Fig. 1.

Our testbed for this problem is a circular maze environment (CME; see Fig. 1), in which the goal is to tip and tilt the maze so as to move a marble from an outer ring into an inner circle. This is an interesting domain for studying real-time control because it is intuitively easy to pick up for people — even children play with similar toys without prior experience with these mazes — and yet is a complex learning domain for artificial agents due to its constrained geometry, underactuated control, nonlinear dynamics, and long planning horizon with several discontinuities [13], [14]. Adding to this challenge, the CME is a system that is usually in motion, so planning and control must be done in real-time, or else the ball will continue to roll in possibly unintended ways.

The learning approach we present in this paper falls under the umbrella of Model-Based Reinforcement Learning (MBRL). In MBRL, a task-agnostic predictive model of the system dynamics is learned from exploration data. This model is then used to synthesize a controller which is used to perform the desired task using a suitable cost function. The model in our case is represented by a physics engine that roughly describes the CME with its physical properties. Additionally, we learn the residual between the actual system and the physics system using Gaussian process regression [9]. Such a combination of a physics engine and a statistical function approximator allows us to efficiently learn models for physical systems while using minimal domain knowledge.

While approaches that combine physical predictions and residuals have been used for control in the past [15], here we demonstrate that this combination can be used as part of a model-predictive controller (MPC) of a much more complex system in real-time. An important point to note here is that the work presented in [15] uses MPC in a discrete action space, whereas for the current system we have to use nonlinear model-predictive control (NMPC) that requires a solution to a nonlinear, continuous control problem in real-time (which requires non-trivial, compute-expensive optimization) [16]. Consequently, the present study deals with a more complicated learning and control problem that is relevant to a wide range of robotic systems. To the best of our knowledge, this is the first time that a hybrid model consisting of a full physics engine along with a statistical function approximator has been used to control a complex physical system in real-time using NMPC. We are also releasing our code for the CME as it is a complex, low-dimensional system that can be used to study real-time physical control<sup>3</sup>.

## II. RELATED WORK

Our work is motivated by the recent advances in (deep) reinforcement learning to solve complex tasks in areas such as computer games [17], [18] and robotics [4], [19]. While these algorithms have been very successful for solving simulated tasks, their applicability in real systems is sometimes questionable due to their relative sample inefficiency. This has

motivated a lot of research in the area of transferring knowledge from a simulation environment to the real world [20], [21], [22], [23], [24], [13]. However, most of these techniques end up being very data intensive, and the agents trained with these algorithms act very differently from how humans solve complex manipulation tasks. Inspired by this mismatch, we attempt to study complex physical puzzles, using model-based agents. Our goal is to understand if a model-based approach is closer to human-like learning.

Recently the robotics community has seen a surge in interest in the use of general-purpose physics engines which can represent complex, multi-body dynamics [25], [26]. These engines have been developed with the intention to allow real-time control of robotic systems while using them as an approximation of the physical world. However, these simulators still cannot model or represent the physical system accurately enough for control, and this has driven a lot of work in the area of sim-to-real transfer [27], [28], [29]. The goal of these methods is to train an agent in simulation and then transfer them to the real system using minimum involvement of the real system during training. However, most of these approaches use a model-free learning approach and thus tend to be sample inefficient. In contrast, we propose a method that trains a MBRL sim-to-real agent and thus achieves very good sample efficiency.

The idea of using residual models for model correction, or hybrid learning models for control of physical systems during learning in physical systems has also been studied in the past [30], [31], [32], [33], [14], [34]. However, most of these studies use prior physics information in the form of differential equations, which requires domain expertise and thus the methods also become very domain specific. While we rely on some amount of domain expertise and assumptions, using a general purpose physics engine to represent the physical system will allow for more readily generalization across a wide range of systems.

A similar CME has been solved with MBRL and deep reinforcement learning, in [14] and [13], respectively. In [14], the analytical equations of motion of the CME have been derived to learn a semi-parametric GP model [35], [36] of the system, and then combined with an optimal controller. In [13], a sim-to-real approach has been proposed, where a policy to control the marble(s) is learned on a simulator from images, and then transferred to the real CME. However, the transfer learning still requires a large amount of data from the real CME.

## III. PROBLEM FORMULATION

We consider the problem of moving the marble to the center of the CME. Our goal is to study the sim-to-real problem as depicted in Fig. 2 in a model-based setting where an agent uses a physics engine as its initial knowledge of the environment’s physics.

Under these settings, we study and attempt to answer the following questions in the present paper.

- 1) What is needed in a model-based sim-to-real architecture for efficient learning in physical systems?

<sup>3</sup><https://www.merl.com/research/license/CME>

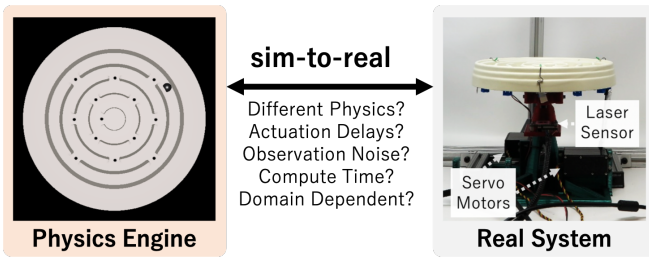


Fig. 2: The sim-to-real problem studied in this paper. We put the real CME on a tip-tilt platform for experiments in the paper. The platform is 3D printed, uses off-the-shelf hobby-grade servo motors and an approximate sensor to measure tip and tilt of the platform. This results in noisy actuations and observations of system states. In addition to different physics (for example, the physics engine cannot accurately model static friction behavior), the simulator has an idealized actuation compared to the real system, i.e., no delays and noise, no approximate tip and tilt sensor.

- 2) How can we design a sim-to-real agent that behaves and learns in a data-efficient, human-like manner?
- 3) How does the performance and learning of our agent compare against the way humans learn to solve these tasks? Can we draw similarities between the human learning policy and the way our RL agent learns?

We use the CME as our test environment for the studies presented in this paper. However, our models and controller design are general-purpose and thus, we expect the proposed techniques could find generalized use in robotic systems. For the rest of the paper, we call the CME together with the tip-tilt platform the circular maze system (CMS).

The goal of the learning agent is to learn an accurate model of the marble dynamics, that can be used in a controller,  $\pi(\mathbf{u}_k|\mathbf{x}_k)$ , in a model-predictive fashion which allows the CMS to choose an action  $\mathbf{u}_k$  given the state observation  $\mathbf{x}_k$  to drive a marble from an initial condition to the target state. We assume that the system is fully defined by the combination of the state  $\mathbf{x}_k$  and the control inputs  $\mathbf{u}_k$ , and it evolves according to the dynamics  $p(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}_k)$  which are composed of the marble dynamics in the maze and the tip-tilt platform dynamics.

As a simplification, we assume that the marble dynamics is independent of the radial dynamics in each of the individual rings, i.e., we quantize the radius of the marble position into the 4 rings of the maze. We include the orientation of the tip-tilt platform as part of the state for our dynamical system, obtaining a five-dimensional state representation for the system, i.e.,  $\mathbf{x} = (r_d, \beta, \gamma, \theta, \dot{\theta})$ . It can be noted that the radius  $r_d$  is a discrete variable, whereas the rest of the state variables are continuous. The terms  $\beta, \gamma$  represent the  $X$  and  $Y$ -orientation of the maze platform, respectively, and  $\theta, \dot{\theta}$  represent the angular position and velocity of the marble, measured with respect to a fixed frame of reference. Since  $r_d$  is fixed for each ring of the CME, we remove  $r_d$  from the state representation of the CMS for the rest of the paper.

Thus, the state is represented by a four-dimensional vector  $\mathbf{x} = (\beta, \gamma, \theta, \dot{\theta})$ . As can be seen in Fig. 2, the angles  $\beta, \gamma$  are measured using a laser sensor that is mounted on the tip-tilt platform while the state of the ball could be observed from a camera mounted above the CMS. For more details, interested readers are referred to [14].

We assume that there is a discrete planner, which can return a sequence of gates that the marble can then follow to move to the center. Furthermore, from the human experiments we have observed that human subjects always try to bring the marble in front of the gate, and then tilt the CME to move it to the next ring. Therefore, we design a lower level controller to move the marble to the next ring when the marble is placed in front of the gate to the next ring. Thus, the task of the learned controller is to move the marble in a controlled way so that it can transition through the sequence of gates to reach the center of the CME.

Before describing our approach, we introduce additional nomenclature we will use in this paper. We represent the physics engine by  $f^{\text{PE}}$ , the residual dynamics model by  $f^{\text{GP}}$ , and the real system model by  $f^{\text{real}}$ , such that  $f^{\text{real}}(\mathbf{x}_k, \mathbf{u}_k) \approx f^{\text{PE}}(\mathbf{x}_k, \mathbf{u}_k) + f^{\text{GP}}(\mathbf{x}_k, \mathbf{u}_k)$ . We use MuJoCo [25] as the physics engine, however, we note that our approach is agnostic to the choice of physics engine. In the following sections, we describe how we design our sim-to-real agent in simulation, as well as on the real system.

#### IV. APPROACH

Our approach for designing the learning agent is inspired by human physical reasoning: people can solve novel manipulation tasks with a handful of trials. This is mainly because we rely on already-learned notions of physics. Following a similar principle, we design an agent whose notion of physics comes from a physics engine. The proposed approach is shown as a schematic in Fig. 3.

As described earlier in Fig 2, we want to design a sim-to-real agent, which can bridge the gap between the simulation environment and the real world in a principled fashion. The gap between the simulated environment and the real world can be attributed to mainly two factors. First, physics engines represent an approximation of the physics of the real systems, because they are designed based on limited laws of physics, domain knowledge, and convenient approximations often made for mathematical tractability. Second, there are additional errors due to system-level problems, such as observation noise and delays, actuation noise and delays, finite computation time to update controllers based on observations, etc.

Consequently, we train our agent to bridge the sim-to-real gap by first estimating the parameters of the physics engine, and then compensate for the different system-level problems as the agent tries to interact with the real system. We use a Gaussian process model to model the residual dynamics of the real system that cannot be modeled by the best estimated parameters of the physics engine. In the rest of this section, we describe the details of the physics engine for the CME,

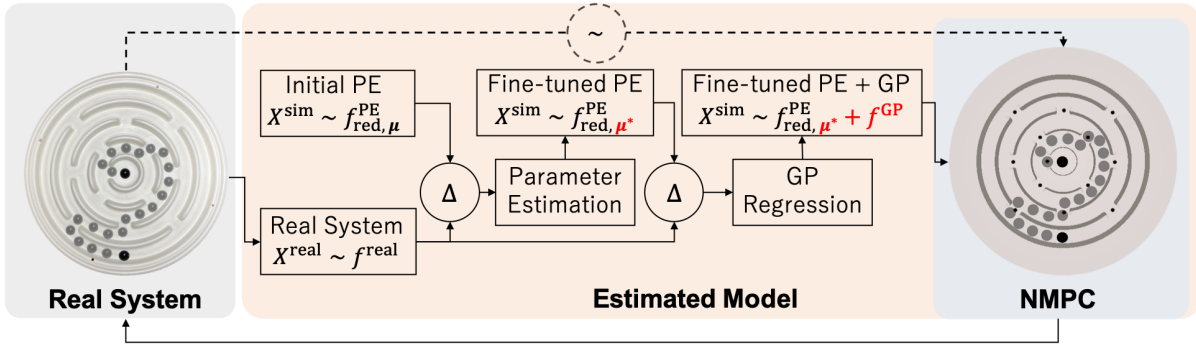


Fig. 3: The learning approach used in this paper to create a predictive model for the physics of the CME in the real system. We create a predictive model for the marble dynamics in the CME using a physics engine. We start with a MuJoCo-based physics engine (PE) with random initial parameters for dynamics, and estimate these parameters  $\mu^*$  from the residual error between simulated and real CME using CMA-ES. The remaining residual error between simulated and real CME is then compensated by using Gaussian process (GP) regression during iterative learning. Finally, we use the estimated model to control the real CME with NMPC policy.

and provide our approach for correcting the physics engine as well as modeling other system-level issues with the CMS.

#### A. Physics Engine Model Description

As described earlier, we use MuJoCo [25] as our physics engine,  $f^{\text{PE}}$ . Note that in our model we ignore the radial movement of the marble in each ring, and describe the state only with the angular position of the marble as described in Sec. III. Consequently, we restrict the physics engine to consider only the angular dynamics of the marble in each ring, i.e., the radius of the marble position is fixed. However, in order to study the performance of the agent in simulation, we also create a full model of the CME where the marble does not have the angular state constraint. Thus, we create two different physics engine models:  $f_{\text{red}}^{\text{PE}}$  represents the reduced physics engine available to our RL model, and  $f_{\text{full}}^{\text{PE}}$  uses the full internal state of the simulator.  $f_{\text{red}}^{\text{PE}}$  differs from  $f_{\text{full}}^{\text{PE}}$  in two key ways. In the forward dynamics of the  $f_{\text{red}}^{\text{PE}}$  model, we set the location of the marble to be in the center of each ring, while this is tracked in  $f_{\text{full}}^{\text{PE}}$ . Additionally, because we cannot observe the spin of the ball in real experiments, we do not include it in  $f_{\text{red}}^{\text{PE}}$ , while it is included in  $f_{\text{full}}^{\text{PE}}$ . We use this  $f_{\text{full}}^{\text{PE}}$  model for analyzing the behavior of our agent in the preliminary studies in simulation. This serves as an analog to the real system in the simulation studies we present in the paper. We call this set of experiments *sim-to-sim*. These experiments are done to determine whether the agent can successfully adapt its physics engine when initialized with an approximation of a more complicated environment.

#### B. Model Learning

We consider a discrete-time system:

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k) + \mathbf{e}_k, \quad (1)$$

where  $\mathbf{x}_k \in \mathbb{R}^4$  denotes the state,  $\mathbf{u}_k \in \mathbb{R}^2$  the actions, and  $\mathbf{e}_k$  is assumed to be a zero mean white Gaussian noise with diagonal covariance that represents the uncertainty about the state at the discrete time instant  $k \in [1, \dots, T]$ .

---

#### Algorithm 1 Model learning procedure

---

- 1: Collect  $N$  episodes in the real system using Alg. 2
  - 2: Compute simulator trajectories as  $f_{\text{red}, \mu}^{\text{PE}}(\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}})$ , from the real system  $N$  episodes
  - 3: *Estimate physical parameters* using CMA-ES
  - 4: **while** Model performance not converged **do**
  - 5:   Collect  $N$  episodes in CMS using Alg. 2
  - 6:   Compute simulator trajectories  $\mathbf{x}_{k+1}^{\text{sim}}$  for data in  $D$
  - 7:   Train residual GP model
  - 8: **end while**
- 

---

#### Algorithm 2 Rollout an episode using NMPC

---

- 1: Initialize time index  $k \leftarrow 0$
  - 2: Reset the real system by randomly placing the marble to outermost ring
  - 3: **while** The marble does not reach innermost ring **and** not exceed time limit **do**
  - 4:   Set real state to simulator  $\mathbf{x}_k^{\text{sim}} \leftarrow \mathbf{x}_k^{\text{real}}$
  - 5:   Compute trajectory  $(X^{\text{sim}}, U^{\text{sim}})$  using NMPC
  - 6:   Apply initial action  $\mathbf{u}_k^{\text{real}} = \mathbf{u}_0^{\text{sim}}$  to the real system
  - 7:   Store transition  $D \leftarrow D \cup \{\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}}, \mathbf{x}_{k+1}^{\text{real}}\}$
  - 8:   Increment time step  $k \leftarrow k + 1$
  - 9: **end while**
- 

In the proposed approach, the unknown dynamics  $f$  in Eq. 1 represents the CMS dynamics,  $f^{\text{real}}$ , and it is modeled as the sum of two components:

$$f^{\text{real}}(\mathbf{x}_k, \mathbf{u}_k) \approx f_{\text{red}}^{\text{PE}}(\mathbf{x}_k, \mathbf{u}_k) + f^{\text{GP}}(\mathbf{x}_k, \mathbf{u}_k), \quad (2)$$

where  $f_{\text{red}}^{\text{PE}}$  denotes the physics engine model defined in the previous section, and  $f^{\text{GP}}$  denotes a Gaussian process model that learns the residual between real dynamics and simulator dynamics. We learn both the components  $f_{\text{red}}^{\text{PE}}$  and  $f^{\text{GP}}$  to improve model accuracy. The approach is presented as pseudo-code in Algorithm 1 and described as follows.

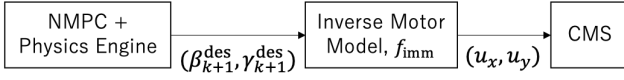


Fig. 4: The agent learns an inverse model of the servo motors to compensate for the delay in actuation when interacting with the real system.

1) *Physical Parameter Estimation:* We first estimate physical parameters of the real system. As measuring physical parameters directly in the real system is difficult, we estimate four friction parameters of MuJoCo by using CMA-ES [37]. More formally, we denote the physical parameters as  $\boldsymbol{\mu} \in \mathbb{R}^4$ , and the physics engine with the parameters as  $f_{\text{red},\boldsymbol{\mu}}^{\text{PE}}$ .

As described in Algorithm 1, we first collect multiple episodes with the real system using the NMPC controller described in Sec. IV-D. Then, CMA-ES is used to estimate the best friction parameters  $\boldsymbol{\mu}^*$  that minimizes the difference between the movement of the marble in the real system and in simulation as:

$$\boldsymbol{\mu}^* = \arg \min_{\boldsymbol{\mu}} \frac{1}{\|D\|} \sum_{(\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}}, \mathbf{x}_{k+1}^{\text{real}}) \in D} \|\mathbf{x}_{k+1}^{\text{real}} - f_{\text{red},\boldsymbol{\mu}}^{\text{PE}}(\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}})\|_{W_{\boldsymbol{\mu}}}^2, \quad (3)$$

where  $D$  represents the collected transitions in the real system,  $W_{\boldsymbol{\mu}}$  is the weight matrix whose value is 1 only related to the angular position term of the marble  $\theta_{k+1}$  in the state  $\mathbf{x}_{k+1}$ .

2) *Residual Model Learning Using Gaussian Process:* After estimating the physical parameters, a mismatch remains between the simulator and the real system because of the modeling limitations described in the beginning of this section. To get a more accurate model, we train a Gaussian Process (GP) model via marginal likelihood maximization [9], with a standard linear kernel, to learn the residual between the two systems by minimizing the following objective:

$$L^{\text{GP}} = \frac{1}{\|D\|} \sum_{(\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}}, \mathbf{x}_{k+1}^{\text{real}}) \in D} \|\mathbf{x}_{k+1}^{\text{real}} - f_{\text{red},\boldsymbol{\mu}^*}^{\text{PE}}(\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}}) - f^{\text{GP}}(\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}})\|^2. \quad (4)$$

Note that after collecting the trajectories in the real system, we collect the simulator estimates of the next state  $\mathbf{x}_{k+1}^{\text{sim}}$  using the physics engine with the estimated physical parameters  $\boldsymbol{\mu}^*$ . This is done by resetting the state of the simulator to every state  $\mathbf{x}_k^{\text{real}}$  along the collected trajectory and applying the action  $\mathbf{u}_k^{\text{real}}$  to obtain the resulted next state  $\mathbf{x}_{k+1}^{\text{sim}} = f_{\text{red},\boldsymbol{\mu}^*}^{\text{PE}}(\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}})$ , and store the tuple  $\{\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}}, \mathbf{x}_{k+1}^{\text{sim}}\}$ . Thus, the GPs learn the input-output relationship:  $f^{\text{GP}}(\mathbf{x}_k^{\text{real}}, \mathbf{u}_k^{\text{real}}) = \mathbf{x}_{k+1}^{\text{real}} - \mathbf{x}_{k+1}^{\text{sim}}$ . Two independent GP models are trained, one each for the position and velocity of the marble.

3) *Modeling Motor Behavior:* As shown in Fig. 2, the tip-tilt platform in the CMS is actuated by hobby-grade servo motors which work in position control mode. These motors use a controller with a finite settling time which is longer

than the control interval used in our experiments. This results in actuation delays for the action computed by any control algorithm, and the platform always has non-zero velocity. The physics engine, on the other hand, works in discrete time and thus the CME comes to a complete rest after completing a given action in a control interval. Consequently, there is a discrepancy between the simulation and the real system in the sense that the real system gets delayed actions. To compensate for this problem, we learn an inverse model for motor actuation. This inverse model of the motor predicts the action to be sent to the motors for the tip-tilt platform to achieve a desired state  $(\beta_{k+1}^{\text{des}}, \gamma_{k+1}^{\text{des}})$  given the current state  $(\beta_k, \gamma_k)$  at instant  $k$ . Thus, the control signals computed by the optimization process are passed through this function that generates the commands  $(u_x, u_y)$  for the servo motors. This is also shown as a schematic in Fig. 4 for clarity. We represent this inverse motor model by  $f_{\text{imm}}$ . The motor model  $f_{\text{imm}}$  is learned using a standard autoregressive model with external input. This is learned by collecting motor response data by exciting the CMS using sinusoidal inputs for the motors before the model learning procedure in Algorithm 1.

### C. Trajectory Optimization using iLQR

We use the iterative LQR (iLQR) as the optimization algorithm for model-based control [38]. While there exist optimization solvers which can generate better optimal solutions for model-based control [39], we use iLQR as it provides a compute-efficient way of solving the optimization problem for designing the controller. Formally, we solve the following *trajectory optimization problem* to manipulate the controls  $\mathbf{u}_k$  over a certain number of time steps  $[T-1]$

$$\begin{aligned} \min_{\mathbf{x}_k, \mathbf{u}_k} \quad & \sum_{k \in [T]} \ell(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t.} \quad & \mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k) \\ & \mathbf{x}_0 = \tilde{\mathbf{x}}_0. \end{aligned} \quad (5)$$

For the state cost, we use a quadratic cost function for the state error measured from the target state  $\mathbf{x}_{\text{target}}$  (which in the current case is the nearest gate for the marble), as represented by the following equation:

$$\ell(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}_{\text{target}}\|_W^2, \quad (6)$$

where the matrix  $W$  represents weights used for different states. For the control cost, we penalize the control using a quadratic cost as well, given by the following equation:

$$\ell(\mathbf{u}) = \lambda_{\mathbf{u}} \|\mathbf{u}\|^2. \quad (7)$$

We tried using different smoother versions of the cost function [38] but it did not change the behavior of the iLQR optimization. The discrete-time dynamics  $\mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k)$  and the cost function are used to compute locally linear models and a quadratic cost function for the system along a trajectory. These linear models are then used to compute optimal control inputs and local gain matrices by iteratively solving the associated LQR problem. For more details of iLQR, interested readers are referred to [38]. The solution

to the trajectory optimization problem returns an optimal sequence of states and control inputs for the system to follow. We call this the reference trajectory for the system, denoted by  $X^{\text{ref}} \equiv \mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T$ , and  $U^{\text{ref}} \equiv \mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{T-1}$ . The matrix  $W$  used for the experiments is diagonal,  $W = \text{diag}(4, 4, 1, 0.4)$  and  $\lambda_u = 20$ .

#### D. Online Control using Nonlinear Model-Predictive Control

While it is easy to control the movement of the marble in the simulation environment, controlling the movement of the marble in the real system is much more challenging. This is mainly due to complications such as static friction (which remains poorly modeled by the physics engine), or delays in actuation. As a result, the real system requires online model-based feedback control. While re-computing an entire new trajectory upon a new observation would be the optimal strategy, due to lack of computation time in the real system, we use a trajectory-tracking MPC controller. We use an iLQR-based NMPC controller to track the trajectory obtained from the trajectory optimization module to control the system in real-time. The controller uses the least-squares tracking cost function given by the following equation:

$$\ell_{\text{tracking}}(\mathbf{x}) = \|\mathbf{x}_k - \mathbf{x}_k^{\text{ref}}\|_Q^2, \quad (8)$$

where  $\mathbf{x}_k$  is the system state at instant  $k$ ,  $\mathbf{x}_k^{\text{ref}}$  is the reference state at instant  $k$ , and the matrix  $Q$  is a weight matrix. The cost on the control actions remain the same as during trajectory optimization. The system trajectory is rolled out forward in time from the observed state, and the objective in Eq. 8 is minimized to obtain the desired control signals.

We implement the control on both the real and the simulation environment at a control rate of 30 Hz. As a result, there is not enough time for the optimizer to converge to the optimal feedback solution. Thus, we warm-start the optimizer with a previously computed trajectory. Furthermore, the derivatives during the system linearization in the backward step of iLQR and the forward rollout of the iLQR are computed using parallel computing in order to compute good solutions within the time provided to compute a feedback step.

## V. EXPERIMENTS

In this section, we try to answer the following questions with our experiments to describe the performance of our agent.

- Can we learn to solve the desired task of moving the marble to the center of the maze within minutes of interacting with the environment?
- How does the proposed method compare against the learning behavior shown by human subjects?

#### A. Physical Property Estimation using CMA-ES

We first demonstrate how physical parameter estimation works in two different environments; sim-to-sim and sim-to-real settings. For sim-to-sim setting, we regard the full model  $f_{\text{full}}^{\text{PE}}$  as a real system because it contains full internal state that is difficult to observe in the real setup as described in

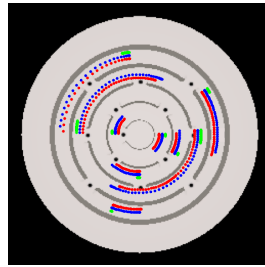


Fig. 5: Comparison of real trajectories (red), predicted trajectories (blue) using the estimated physical properties using CMA-ES, and trajectories using the default physical properties (green) in the sim-to-sim experiment. The trajectories are generated with a random policy from random initial points.

Sec. IV-A. Also, we regard the reduced model  $f_{\text{red}}^{\text{PE}}$ , which has the same state that can be observed in the real system, as a simulator. For sim-to-real setting, we measure the difference between the real system and the reduced model  $f_{\text{red}}^{\text{PE}}$ . For  $f_{\text{red}}^{\text{PE}}$ , we start with default values given by MuJoCo, and we set smaller friction parameters to  $f_{\text{full}}^{\text{PE}}$  in the sim-to-sim setting, because we found the real maze board is much more slippery than what default MuJoCo's parameters would imply.

We collected samples using the NMPC controller computed using current  $f_{\text{red}}^{\text{PE}}$  models, and found the objective defined in equation 3 converges only  $\sim 10$  transitions for each ring. For sim-to-sim experiment, the RMSE of ball location  $\theta$  in two dynamics becomes  $\approx 2e - 3$  [rad] ( $\approx 0.1$  [deg]), which we conclude the CMA-ES produces accurate enough parameters. Figure 5 shows the real trajectories obtained by  $f_{\text{full}}^{\text{PE}}$  (in red), simulated trajectories obtained by  $f_{\text{red}}^{\text{PE}}$  with optimized friction parameters (in blue), and simulated trajectories before estimating friction parameters (in green). This qualitatively shows that the estimated friction parameters successfully bridge the gap between two different dynamics. Since tuning friction parameters for MuJoCo is not intuitive, it is evident that we can rely on CMA-ES to determine more optimal friction parameters instead. Similarly, we find that sim-to-real experiment, the RMSE of ball position  $\theta$  between the physics engine and real system decreased to  $\approx 9e - 3$  [rad] after CMA-ES optimization. However, we believe this error still diverges in rollout and we still suffer from static friction. We also observed that CMA-ES optimization in the sim-to-real experiments quickly finds a local minima with very few samples, and further warm starting the optimization with more data results in another set of parameters for the physics engine with similar discrepancy between the physics engine and the real system. Thus, we perform the CMA-ES parameter estimation only once in the beginning and more finetuning to GP regression.

#### B. Control Performance on Real System

We found the *sim-to-sim* agent learns to perform well with just CMA-ES finetuning, and thus we skip further control results for the *sim-to-sim* agent, and only present results on the real system with additional residual learning for improved performance. While CMA-ES works well in the *sim-to-sim*

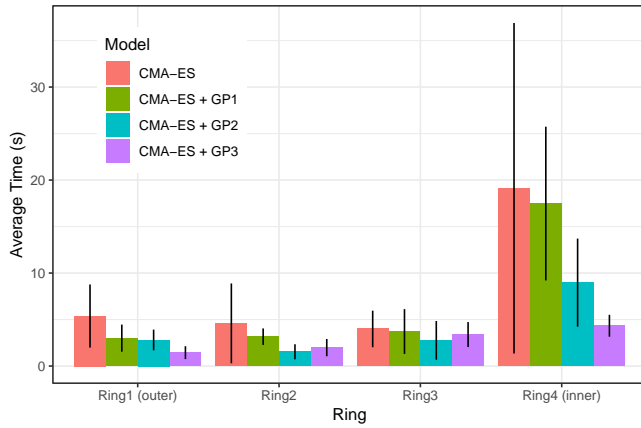


Fig. 6: Comparison of average time spent by the marble in each ring during learning. This plot shows the improvement in the performance of the controller upon learning of the residual model. Note that the controller completely fails without CMA-ES initialization, and thus, those results are not included.

transfer problem, if we want a robot to solve the CME, there will necessarily be differences between the internal model and real-world dynamics. We take inspiration from how people understand dynamics – they can both capture physical properties of items in the world, and also learn the dynamics of arbitrary objects and scenes. For this reason we augmented the CMA-ES model with machine learning data-driven models that can improve the model accuracy as more experience (data) is acquired. We opted for GP as data-driven models because of their high flexibility in describing data distribution and data efficiency [40].

We considered models with incrementally more learning for the data-driven residual model on top of the base CMA-ES, in batches of 5 attempts. Namely, ‘CMA-ES’ represents the base CMA-ES model without any residual modeling, while ‘CMA-ES + GP1’ represents a model that has learned a residual model from watching 5 attempts of the ‘CMA-ES’ model. Similarly, ‘CMA-ES + GP2’ and ‘CMA-ES + GP3’ learn the residual distribution from 10 experiments (5 with ‘CMA-ES’ and 5 with ‘CMA-ES + GP1’) and 15 experiments (5 each from ‘CMA-ES’, ‘CMA-ES + GP1’ and ‘CMA-ES + GP2’), respectively.

Figure 6 shows the average time spent in each ring for all four levels of training. As expected, models trained with a larger amount of data consistently improve the performance, i.e., spending less time in each ring. The improvement in performance can be seen especially in the outermost (Ring1) and innermost ring (Ring4). The outermost ring has the largest radius and is more prone to oscillations, which the model learns to control. Similarly, in the innermost ring, static friction causes small actions to have larger effects. In the middle two rings, where the ‘CMA-ES’ was already performing well, we can still observe an improvement in the reduction of the variability of the performance. A video showing the performance at different stages of learning could

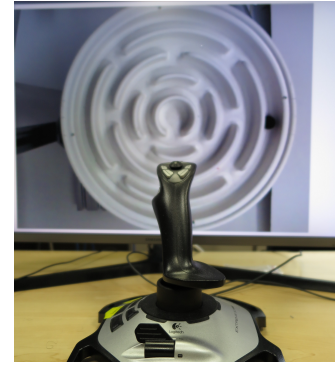


Fig. 7: The joystick system along with video feed of the CME used for human experiments. Human subjects were asked to solve the maze by looking at the video feed of the marble movement. Similar to the MPC policy, the joystick allows control of the two motors of the CMS, and is very natural to control for lots of humans.

be found here<sup>4</sup>.

### C. Comparison with Human Performance

To compare our system’s performance against human learning, we asked 15 participants to perform a similar CME task. Since the learning algorithm uses the two servo motors to control the motion of the marble, the human subjects were asked to control the CMS with a 2 DoF joystick (see Fig. 7). To familiarize participants with the controls, they were given one minute to play with the maze using the joystick, but there was no ball in the maze during this time. This is similar to how the model pre-learned the inverse motor model  $f_{imm}$  without learning ball dynamics. Afterwards, the ball was placed at a random point in the outermost ring, and participants were asked to guide the ball to the center of the maze. They were asked to solve the CME five times, and we recorded how long they took for each solution and how much time the ball spent in each ring. Two participants were excluded from analysis because they could not solve the maze five times within the 15 minutes allotted to them.

Because people were given five maze attempts (and thus between zero and four prior chances to learn during each attempt), we compare human performance against the CMA-ES and CMA-ES+GP1 versions of our model that have comparable amounts of training.

We find that while there was a slight numerical decrease in participants’ solution times over the course of the five trials, this did not reach statistical reliability ( $\chi^2(1) = 1.63$ ,  $p = 0.2$ ): participants spent an average of 110 seconds (95%CI : [66, 153]) to solve the maze the first time, and 79 seconds (95%CI : [38, 120]) to solve the maze the last time, and only 8 of 13 participants solved the maze faster on the last trial as compared to their first. This is similar to the learning pattern found in our model, where the solution time decreased from 33s using CMA-ES to 27s using CMA-ES+GP1, which similarly was not statistically reliable ( $t(15) = 0.56$ ,  $p =$

<sup>4</sup><https://youtu.be/xaxNCXBovpc>

TABLE I: Average time spent in each ring [sec].

	Human	CMA-ES + GP0/1
Ring 1 (outermost ring)	22.6	4.18
Ring 2	8.0	3.87
Ring 3	24.3	3.85
Ring 4 (innermost ring)	41.1	18.29

0.58). A qualitative comparison between the human subjects and our proposed method could be seen in the videos<sup>5</sup>.

In addition, Table I shows the time that people and the model kept the ball in each ring. For statistical power we have averaged over all human attempts, and across CMA-ES and CMA-ES + GP1 to equate to human learning. In debriefing interviews, participants indicated that they found that solving the innermost ring was the most difficult, as indicated by spending more time in that ring than any others (all  $ps < 0.05$  by Tukey HSD pairwise comparisons). This is likely because small movements will have the largest effect on the marble’s radial position, requiring precise prediction and control. Similar to people, the model also spends the most time in the inner ring (all  $ps < 0.002$  by Tukey HSD pairwise comparisons), suggesting that it shares the same prediction and control challenges as people. In contrast, a fully trained standard reinforcement learning algorithm – the soft actor-critic (SAC) [41] – learns a different type of control policy in simulation and spends the *least* amount of time in the innermost ring, since the marble has the shortest distance to travel (see Appendix A for more detail).

## VI. CONCLUSIONS AND FUTURE WORK

Humans can learn and adapt their approach to perform complex tasks within minutes of interaction with a novel system. Studies from cognitive science suggest that this is because people have internal models of physics that are well calibrated to the world [12]. There has been much recent interest in using this idea in robotic systems by basing planning around physics engines. This is mainly based on the vision that these physics engines can be used for real-time control of robotic systems by providing the capability for real-time physical reasoning.

Here we use this idea to design agents that can interact with the world in a human-like fashion. We presented a learning method for navigating a marble in a complex circular maze environment. Learning consists of initializing a physics engine, where the physics parameters are initially estimated using the real system. The error in the physics engine is then compensated using a Gaussian process regression model which is used to model the residual dynamics. These models are used to control the marble in the maze environment using iLQR in a feedback MPC fashion. We showed that the proposed method can learn to solve the task of driving the marble to the center of the maze within a few minutes of interacting with the system. We contrasted the learning behavior against the time taken by humans to solve the problem to show comparable behavior.

<sup>5</sup><https://youtu.be/xaxNCXBovpc>

TABLE II: Average time spent each ring in simulation [sec].

	CMA-ES	SAC
Ring 1 (outermost ring)	1.50	0.78
Ring 2	1.00	0.83
Ring 3	2.60	0.86
Ring 4 (innermost ring)	7.17	0.73

One of the benefits of our approach is its flexibility: because it learns based off of a general-purpose physics engine, this approach should generalize well to other real-time physical control tasks. Furthermore, the separation of the dynamics and control policy should facilitate transfer learning. If the maze material or ball were changed (e.g., replacing it with a small die or coin), then the physical properties and residual model would need to be quickly relearned, but the control policy should be relatively similar. In future work, we plan to test the generality and transfer of this approach to different mazes and marbles. For more effective use of physics engines for these kind of problems, we would like to interface general-purpose robotics optimization software [42] to make it more useful for general-purpose robotics application.

## APPENDIX

### A. Control Performance on Simulation

In order to compare the performance of our approach and a model-free RL algorithm, we train a SAC [41] agent with  $f_{\text{full}}^{\text{PE}}$  dynamics in simulation. The hyperparameters, architectures, activation function of SAC are the same as used in [41]. We also evaluate the performance of our method in sim-to-sim setting, which omits the GP part because CMA-ES quickly matches the behavior of the simulator in the sim-to-sim setting, as described in Sec. V-A.<sup>6</sup>

Table. II shows the average time spent in each ring for both methods. The SAC model solves the maze faster than the CMA-ES algorithm, but does so by speeding the ball through each ring in approximately equal time, unlike both CMA-ES and people. This is likely because the SAC agent had extensive experience to learn its control policy: it was trained for five million steps on the simulator, which is equivalent to approximately two days training time if done on a real system.

### B. MuJoCo Model Setting

As written in Sec. IV-A, we prepare two different physics engine models:  $f_{\text{red}}^{\text{PE}}$  and  $f_{\text{full}}^{\text{PE}}$ . Table. III summarizes the friction parameters  $\mu$  used for each environment. We note that these initial parameters are optimized by CMA-ES. We set the same friction parameters to all objects in the simulator: the walls and bottom that construct the circular maze, and the marble. We have modeled the mass and size of the marble, and geometry of the circular maze based on our measurements of the real CME used in CMS.

<sup>6</sup>We attempted to train SAC on the real CME, but were unable to demonstrate any learning after three days, perhaps due to complications like the continuous action space or high control frequency. However, [13] demonstrated sim-to-real with transfer learning could solve a somewhat different CME, suggesting a possible additional comparison for future work.



TABLE III: Physical parameters used in sim-to-sim experiments. The  $f_{\text{red}}^{\text{PE}}$  uses default parameters of MuJoCo, whereas the  $f_{\text{full}}^{\text{PE}}$  is more slippery, because we found that the real model is actually more slippery than what default parameters would imply [43].

	$f_{\text{full}}^{\text{PE}}$	$f_{\text{red}}^{\text{PE}}$
Slide friction	$1e-3$	1
Spin friction	$1e-6$	$5e-3$
Roll friction	$1e-7$	$1e-4$
Friction loss	$1e-6$	0

## REFERENCES

- [1] François Osiurak and Dietmar Heinke. Looking for intoelligence: A unified framework for the cognitive study of human tool use and technology. *American Psychologist*, 73(2):169–185, 2018.
- [2] Kuan Fang, Yuke Zhu, Animesh Garg, Andrey Kurenkov, Viraj Mehta, Li Fei-Fei, and Silvio Savarese. Learning task-oriented grasping for tool manipulation from simulated self-supervision. *The International Journal of Robotics Research*, 2019.
- [3] Marc Toussaint, Kelsey R. Allen, Kevin A. Smith, and Joshua B. Tenenbaum. Differentiable Physics and Stable Modes for Tool-Use and Manipulation Planning. In *Robotics: Science and Systems XIV*. Robotics: Science and Systems Foundation, 2018.
- [4] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- [5] Alvaro Sanchez-Gonzalez, Nicolas Heess, Jost Tobias Springenberg, Josh Merel, Martin Riedmiller, Raia Hadsell, and Peter Battaglia. Graph networks as learnable physics engines for inference and control. volume 80 of *Proceedings of Machine Learning Research*, pages 4470–4479, Stockholm, Sweden, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- [6] Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Józefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, Jonas Schneider, Szymon Sidor, Josh Tobin, Peter Welinder, Lilian Weng, and Wojciech Zaremba. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020.
- [7] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40, 2017.
- [8] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- [9] Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*. MIT press Cambridge, MA, 2006.
- [10] Peter W. Battaglia, Jessica B. Hamrick, and Joshua B. Tenenbaum. Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences*, 110(45):18327–18332, 2013.
- [11] Kevin A. Smith, Peter W. Battaglia, and Edward Vul. Different Physical Intuitions Exist Between Tasks, Not Domains. *Computational Brain & Behavior*, 1(2):101–118, 2018.
- [12] Kelsey R Allen, Kevin A Smith, and Joshua B Tenenbaum. The tools challenge: Rapid trial-and-error learning in physical problem solving. *arXiv preprint arXiv:1907.09620*, 2019.
- [13] J. v. Baar, A. Sullivan, R. Corcodel, D. Jha, D. Romeres, and D. Nikovski. Sim-to-real transfer learning using robustified controllers in robotic tasks involving complex dynamics. In *2019 International Conference on Robotics and Automation (ICRA)*, May 2019.
- [14] D. Romeres, D. K. Jha, A. DallaLibera, B. Yerazunis, and D. Nikovski. Semiparametrical gaussian processes learning of forward dynamical models for navigating in a circular maze. In *2019 International Conference on Robotics and Automation (ICRA)*, May 2019.
- [15] Anurag Ajay, Maria Bauza, Jiajun Wu, Nima Fazeli, Joshua B Tenenbaum, Alberto Rodriguez, and Leslie P Kaelbling. Combining physical simulators and object-based networks for control. *arXiv preprint arXiv:1904.06580*, 2019.
- [16] Moritz Diehl, Hans Joachim Ferreau, and Niels Haverbeke. Efficient numerical methods for nonlinear mpc and moving horizon estimation. In *Nonlinear model predictive control*, pages 391–417. Springer, 2009.
- [17] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [18] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.
- [19] John Schulman, Sergey Levine, Pieter Abbeel, Michael I Jordan, and Philipp Moritz. Trust region policy optimization. In *Icml*, volume 37, pages 1889–1897, 2015.
- [20] Stephen James, Andrew J. Davison, and Edward Johns. Transferring end-to-end visuomotor control from simulation to real world for a multi-stage task. volume 78 of *Proceedings of Machine Learning Research*, pages 334–343. PMLR, 13–15 Nov 2017.
- [21] Florian Golemo, Adrien Ali Taiga, Aaron Courville, and Pierre-Yves Oudeyer. Sim-to-real transfer with neural-augmented robot simulation. In *Conference on Robot Learning*, pages 817–828, 2018.
- [22] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1–8. IEEE, 2018.
- [23] Jie Tan, Tingnan Zhang, Erwin Coumans, Atil Iscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and Vincent Vanhoucke. Sim-to-real: Learning agile locomotion for quadruped robots. In *Robotics: Science and Systems (RSS)*, 2018.
- [24] Y. Chebotar, A. Handa, V. Makoviychuk, M. Macklin, J. Issac, N. Ratliff, and D. Fox. Closing the sim-to-real loop: Adapting simulation randomization with real world experience. In *International Conference on Robotics and Automation (ICRA)*. IEEE, 2019.
- [25] 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*, pages 5026–5033, Oct 2012.
- [26] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. <http://pybullet.org>, 2016–2019.
- [27] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 23–30, 2017.
- [28] Stephen James, Paul Wohlhart, Mrinal Kalakrishnan, Dmitry Kalashnikov, Alex Irpan, Julian Ibarz, Sergey Levine, Raia Hadsell, and Konstantinos Bousmalis. Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12627–12637, 2019.
- [29] Fabio Ramos, Rafael Possas, and Dieter Fox. Bayessim: Adaptive domain randomization via probabilistic inference for robotics simulators. In *Robotics: Science and Systems*, 2019.
- [30] Lukas Hewing, Juraj Kabzan, and Melanie N Zeilinger. Cautious model predictive control using gaussian process regression. *IEEE Transactions on Control Systems Technology*, 2019.
- [31] Matteo Saveriano, Yuchao Yin, Pietro Falco, and Dongheui Lee. Data-efficient control policy search using residual dynamics learning. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4709–4715. IEEE, 2017.
- [32] Anurag Ajay, Jiajun Wu, Nima Fazeli, Maria Bauza, Leslie P Kaelbling, Joshua B Tenenbaum, and Alberto Rodriguez. Augmenting physical simulators with stochastic neural networks: Case study of planar pushing and bouncing. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3066–3073. IEEE, 2018.
- [33] A. Ajay, J. Wu, N. Fazeli, M. Bauza, L. P. Kaelbling, J. B. Tenenbaum, and A. Rodriguez. Augmenting physical simulators with stochastic neural networks: Case study of planar pushing and bouncing. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems*.
- [34] Tingfan Wu and Javier Movellan. Semi-parametric gaussian process for robot system identification. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 725–731. IEEE, 2012.
- [35] D. Romeres, M. Zorzi, R. Camoriano, and A. Chiuso. Online semi-parametric learning for inverse dynamics modeling. In *IEEE 55th Conference on Decision and Control (CDC)*, pages 2945–2950, 2016.
- [36] D. Nguyen-Tuong and J. Peters. Using model knowledge for learning inverse dynamics. In *2010 IEEE International Conference on Robotics and Automation*, pages 2677–2682, 2010.

- [37] Nikolaus Hansen. The cma evolution strategy: a comparing review. In *Towards a new evolutionary computation*. Springer, 2006.
- [38] Yuval Tassa, Tom Erez, and Emanuel Todorov. Synthesis and stabilization of complex behaviors through online trajectory optimization. In *Intelligent Robots and Systems (IROS)*. IEEE, 2012.
- [39] John T Betts. Survey of numerical methods for trajectory optimization. *Journal of guidance, control, and dynamics*, 21(2):193–207, 1998.
- [40] A. Dalla Libera, D. Romeres, D. K. Jha, B. Yezunian, and D. Nikovski. Model-based reinforcement learning for physical systems without velocity and acceleration measurements. *IEEE Robotics and Automation Letters*, 5(2):3548–3555, 2020.
- [41] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey 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*, pages 1861–1870. PMLR, 2018.
- [42] Russ Tedrake and the Drake Development Team. Drake: Model-based design and verification for robotics, 2019.
- [43] Mujoco. <http://www.mujoco.org/>. Accessed: 2020-01-31.