

Preparing for the Unknown: Learning a Universal Policy with Online System Identification

Wenhao Yu¹, Jie Tan², C. Karen Liu¹, and Greg Turk¹

wenhaoyu@gatech.edu, jietan@google.com, karenliu@cc.gatech.edu, turk@cc.gatech.edu

¹Interactive Computing, Georgia Institute of Technology, USA

²Google Brain, Google, USA

Abstract—We present a new method of learning control policies that successfully operate under unknown dynamic models. We create such policies by leveraging a large number of training examples that are generated using a physical simulator. Our system is made of two components: a Universal Policy (UP) and a function for Online System Identification (OSI). We describe our control policy as universal because it is trained over a wide array of dynamic models. These variations in the dynamic model may include differences in mass and inertia of the robots components, variable friction coefficients, or unknown mass of an object to be manipulated. By training the Universal Policy with this variation, the control policy is prepared for a wider array of possible conditions when executed in an unknown environment. The second part of our system uses the recent state and action history of the system to predict the dynamics model parameters μ . The value of μ from the Online System Identification is then provided as input to the control policy (along with the system state). Together, UP-OSI is a robust control policy that can be used across a wide range of dynamic models, and that is also responsive to sudden changes in the environment. We have evaluated the performance of this system on a variety of tasks, including the problem of cart-pole swing-up, the double inverted pendulum, locomotion of a hopper, and block-throwing of a manipulator. UP-OSI is effective at these tasks across a wide range of dynamic models. Moreover, when tested with dynamic models outside of the training range, UP-OSI outperforms the Universal Policy alone, even when UP is given the actual value of the model dynamics. In addition to the benefits of creating more robust controllers, UP-OSI also holds out promise of narrowing the Reality Gap between simulated and real physical systems.

I. INTRODUCTION

Numerical simulation of physical phenomena is a powerful tool that has been embraced by researchers and practitioners in computer animation. The success in simulating highly dynamic motion in computer animation, however, has not been transferred in full to robotics. The discrepancy between what can be achieved in simulation and that in real world is referred as the “Reality Gap” in the Evolutionary Robotics community [22, 14]. Researchers have put forth a long list of possible factors that give rise to the Reality Gap, such as simplified dynamic models, inaccurate model parameters, approximated hardware limitations, the absence of uncertainty and latency in sensors and actuators, and other unmodelled factors. Closing the Reality Gap has recently been a major focus in robotics because the ability to transfer knowledge learned in simulation to the real world can potentially unlock the full capability of deep reinforcement learning for robotic applications.

There are two general approaches to learning control policies for real-world operations. One can design a robust control policy to handle a reasonable amount of noise based on a hypothesized dynamic model and hope that the control policy will succeed in the real world without any intervention. This approach is effective if the hypothesized dynamic model is not too far off from the real one. Unfortunately, this is typically not the case for complex robotic system that are performing dynamic tasks involving contacts. Alternatively, one can learn a more accurate dynamic model from real-world data by alternating between control policy optimization, data collection, and dynamic model fitting (or *system identification*). One potential limiting factor of this approach is the need for a large amount of real-world trials which can be expensive, time-consuming, and sometimes dangerous to the robot, humans and objects in its surroundings.

This paper introduces a different approach to combat the model discrepancy. We present a new algorithm that only uses data from a generic physical simulator to learn motor control policies that successfully operate under unknown dynamic models. Our key idea is to aggressively explore the “virtual world” through physical simulation and precompute many of the possible situations the robot might encounter when operating in real world. Although it is easy to dismiss such an approach based on the curse of dimensionality argument, in light of recent disruptive development in deep learning and the availability of large-scale computing capability, we believe it is time to revisit this pre-computation approach. Suppose that we have a way to precompute the optimal control policy for every dynamic model in a parameterized space. Suppose further that we have a fast method that tells us which dynamic model best fits an observed motion sequence. With these two scenarios, during online execution in the real world, we only need to select the right dynamic model by looking at the recent history of robot motion, and then select the corresponding control policy for that dynamic model to achieve optimal motion. While it might not be possible to pre-learn a control policy or pre-explore the entire space of motions for every possible dynamic model, it might be possible to do so for everything that is important for the task of interest.

To this end, we introduce a new approach that exploits a massive amount of simulated data to learn a) a universal control policy (UP) that is capable under a parameterized

space of dynamic models, when provided with the appropriate dynamic model parameters, and b) an on-line system identification model (OSI) that predicts the dynamic model parameters given the current state and the recent history of state-action pairs. Figure 1 shows the relationship between the universal policy and the online system identification. Once trained, the combined algorithm, UP-OSI, can be executed in an unknown dynamic environment. At each time instance, the dynamic model parameters μ is first predicted by the learned system identification model. The universal control policy then takes the predicted model parameters along with the current state to compute the optimal action (Figure 1).

UP-OSI is sample-efficient by design because the algorithm does not require real-world samples during offline training. Another important advantage of UP-OSI is that it does not require the model parameters to be identified prior to execution. While some model parameters might not change over time (e.g. mass, length of a body part) and can be identified offline, other parameters related to the unknown environment, such as the friction coefficient of the floor or the mass of objects being manipulated, cannot be easily identified in advance. Part of the power of UP-OSI is that it can dynamically adapt to changing factors in the environment.

We evaluate our method by learning dynamic motor skills and executing them under unknown dynamic models in simulation. In each of the examples, the control policy can successfully execute the task without knowing some crucial parameters of the dynamic model, such as the inertial and geometric parameters of the robot, variable friction coefficients in the environment, and other task-related parameters. Furthermore, we demonstrate that UP-OSI can operate successfully outside the space of dynamic models used for training, as well as under sudden changes in the environment.

II. RELATED WORK

A. Deep Reinforcement Learning

In recent years, researchers have used deep reinforcement learning to train highly dynamic motor skills in simulated environments that have high-dimensional state and action spaces [23, 31, 30, 18, 20, 10]. For example, Schulman *et al.* demonstrated learning of fullbody humanoid running and getting up with just feedback from the reward function using Trust Region Policy Optimization (TRPO) [30] and Generalized Advantage Estimation (GAE) [31]. Lillicrap *et al.* [18] extended their work of Deep Q-Learning [19] and Deterministic Policy Gradient (DPG) [32] to learn robotic motor skills such as hopping, reaching and 2D walking directly from pixel input. These methods usually require a large amount of interaction time between the agent and the environment, which poses a significant challenge to the robot as well as the experimenter when applying them directly to learn a real-world robotic control task.

In addition, progress has been made in directly learning neural network control policies of manipulation tasks for real robots [16, 25, 17]. While the results are impressive, these methods usually require extensive amount of experimental data

[17, 25] or relatively restrictive settings [16]. It is unclear whether these method would work directly on more dynamic motor skills in the real-world, such as locomotion.

B. Transfer Learning in Reinforcement Learning

Transferring policy learned in simulation to real-world robot has the potential to address the problem of learning complex motor skills for real robots. Much previous effort focused on classical *system identification*, which provides a framework to address the general problem of model inconsistency. In practice, system identification is often interleaved with control policy optimization to minimize the number of required real-world experiments [3, 9, 5]. Some widely-used models, such as linear models [3], Gaussian processes [8, 11], and differential equations [35, 4], have proven effective for continuous dynamics and control tasks with relatively low action space. For example, Abbeel and Ng [3] used a time-variant linear function to model the dynamics of a helicopter from the real-world data while learning a control policy to perform inverted autonomous helicopter flight. Deisenroth and Rasmussen [8] trained a Gaussian process model from real-world data to analytically calculate the control policy gradient, significantly reducing the number of samples compared to sampling-based policy gradient estimation. Ross and Bagnell [28] provided a proof that such iterative processes can converge to an optimal policy, given an accurate dynamic model learning method and a good policy search algorithm. More recently, deep neural networks have been applied to learn both forward dynamics [26] and inverse dynamics [6] from the real-world data, which can potentially model more complicated dynamics. A key drawback of these methods is that their success depends on the quantity and the quality of the real-world data. For highly dynamic or contact-rich tasks, learning an accurate dynamic model and control policy would require a large amount of high quality data, which can be difficult to acquire.

Another line of research utilizes simulated data to train a policy and then directly applies it or adapts it to the real world. James and Johns [13] demonstrated a simulation-trained manipulation controller learned from photo-realistic rendering data that show similar behaviors in the real-world. Rusu *et al.* [29] used progressive networks to efficiently learn a manipulation task on a Jaco arm from a policy that was trained in simulation. One important assumption made by this approach is that the dynamics modeled by the simulator is similar to the real world dynamics. It is unclear whether these methods can work in the situations where the policy is sensitive to the discrepancy in the assumed dynamics model.

In this work, we avoid the need of explicit modeling of system dynamics by learning a Universal Policy (UP) that simultaneously optimizes a wide range of model parameters and an Online System Identification (OSI) network that estimates the dynamic parameters during execution. The concept of UP is similar to Mordatch *et al.* [21] who optimized the motion trajectory for an ensemble of dynamic models perturbed from the assumed one. Their method shows high success rate when tracking a real-world reference trajectory. The motor skills we

are studying require a complex feedback policy, which cannot be achieved by tracking a reference trajectory. Rajeswaran et al. [27] also proposed training a control policy using an ensemble of dynamic models. They used an adversarial training scheme to improve the policy performance while iteratively update the distribution of the dynamic models using data from the target environment. The resulting controller is robust to unknown parameters. In contrast, training UP to achieve a high reward for a variety of dynamic models is an easier learning task because UP explicitly takes model parameters as part of the input.

C. Learning Policy in Unknown Environment

Our work is also related to learning a control policy in an environment with unknown parameters, i.e. Partial Observable Markov Decision Process (POMDP). One example of such work is the idea of event-learning proposed by Szita et al. [34]. They demonstrated that when combined with Static and Dynamic State (SDS) controller, their method performs better than standard method like SARSA in an dynamically varying environment. In contrast to their work, we explicitly incorporate the model parameters as the input to the control policy. With this additional information as input, the control policy has the potential to achieve better performance for a larger range of model parameters. By combining L_1 adaptive control with Iterative Learning Control framework, Pereira et al. [24] demonstrated learning a robust trajectory tracking controller under unknown dynamics and disturbances. However, it is unclear whether this method can be extended to other control tasks such as locomotion.

Based on the Deterministic Policy Gradient method [32], Heess et al. [12] proposed a method that represents a control policy as a recurrent neural network. They demonstrated control policies that identify and memorize task-related signals presented in the observable data sequences, such as the mass of the robot or the position of the target to be grasped. Our work shares a similar goal in learning controllers that can be applied to systems with unknown dynamic model parameters. However, their focus is on learning memory-related tasks rather than transfer to unseen dynamic environments. On the other hand, our work currently does not handle memory-related tasks as our controller only takes a small window of motion history as input.

III. METHODS

Our algorithm consists of two components, a universal policy (UP) and an on-line system identification model (OSI), shown in Figure 1. First, we formulate a reinforcement learning problem to learn a universal policy (UP), $\pi : (\mathbf{x}, \boldsymbol{\mu}) \mapsto \mathbf{u}$, for a space of dynamic models, $\mathbf{x}_{t+1} = f_{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{u}_t)$, parameterized by the dynamic model parameters $\boldsymbol{\mu}$. Unlike a conventional control policy which maps a state vector \mathbf{x} to a control vector \mathbf{u} , UP takes as input both the state and the dynamic model parameters (i.e. $\boldsymbol{\mu}$) and outputs a control vector. Second, we formulate a supervised learning problem to train an online system identification model (OSI),

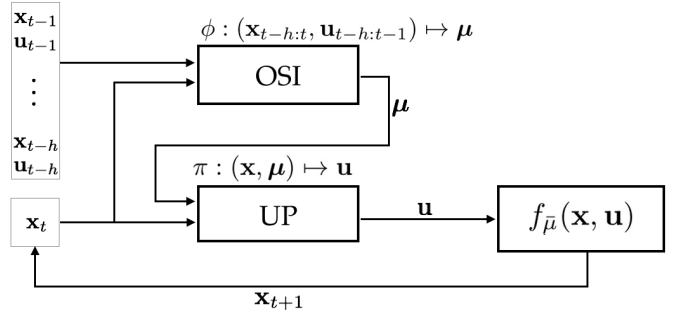


Fig. 1. Overview of UP-OSI. The online system identification model (OSI) takes as input the recent history of the motion and identify the model parameters $\boldsymbol{\mu}$. The universal control policy (UP) then takes the predicted model parameters along with the current state \mathbf{x} to compute the optimal control \mathbf{u} .

$\phi : (\mathbf{x}_{t-h:t}, \mathbf{u}_{t-h:t-1}) \mapsto \boldsymbol{\mu}$, that predicts the dynamic model parameters $\boldsymbol{\mu}$, given the current state \mathbf{x}_t and the past h time instances of the state-action pairs. Unlike the conventional system identification approach, our goal is not to identify a particular system, but to create a *function* that can identify the model parameters for any given trajectory. Both components are represented as a standard neural network and trained offline using simulated data only.

Putting UP and OSI together, at every time instance, we first use OSI (ϕ) to predict the dynamic model parameters $\boldsymbol{\mu}$ based on the current state of the robot \mathbf{x}_t and the recent history of motion $(\mathbf{x}_{t-1}, \mathbf{u}_{t-1}, \dots, \mathbf{x}_{t-h}, \mathbf{u}_{t-h})$. Once $\boldsymbol{\mu}$ is identified, we feed both $\boldsymbol{\mu}$ and the current state \mathbf{x}_t into UP (π) to evaluate the the optimal action \mathbf{u}_t under the predicted dynamic model. We execute \mathbf{u}_t on the robot and push \mathbf{x}_t and \mathbf{u}_t into the history queue. The new state of the system becomes the current state \mathbf{x}_t and the algorithm advances to the next time step (Fig 1).

We propose a novel framework to train the control policy and the system identification model. Conventional approaches alternate between system identification and control policy optimization, which requires a reasonable initial policy or/and an initial system identification model to optimize efficiently. Our method breaks this mutual dependence by first training UP preemptively to anticipate all of the possible dynamic models that OSI might explore during optimization.

A. Learning Universal Policy

Our goal is to learn a control policy that can be generalized to a parameterized space of dynamic models. Many existing methods [7, 33] employ an ensemble approach by learning a discrete set of control policies and consolidating them into one regression model. Our initial attempt with the ensemble approach showed that, for many dynamic tasks, sometimes a small change in the model parameter requires a drastically different control policy to succeed at the given task. Fitting a regression model to this non-smooth landscape of policies often yields poor generalizability.

In this work, we found that it is possible to directly train a large neural net to represent a universal control policy, $\pi(\mathbf{x}, \boldsymbol{\mu})$, for a space of dynamic models parameterized by $\boldsymbol{\mu}$. With a powerful policy optimization algorithm and sufficient data, the

universal policy can achieve high rewards across the space of μ , with comparable performance to policies that have been trained for a specific μ .

We use the Trust Region Policy Optimization (TRPO) method [30] and show that by simply appending the model parameters μ to the input state, TRPO can successfully train a universal control policy. However, we need to modify the exploration scheme of TRPO because the part of the state space that represents μ is not affected by forward simulation when generating rollouts. For each rollout, our algorithm (Algorithm 1) samples μ_i from a uniform distribution ρ_μ , and generate the motion sequence under the policy $\pi(\mathbf{x}, \mu_i)$ and the dynamic model f_{μ_i} . Once the state-action pairs are collected in this manner, the update of π follows TRPO exactly.

B. Learning Online System Identification Model

Even with the ability to perform control under different dynamic models, UP can only succeed at a task when given accurate model parameters, and this information is typically not readily available. We propose to learn an online system identification model (OSI), $\phi : (\mathbf{x}_{t-h:t}, \mathbf{u}_{t-h:t-1}) \mapsto \mu$, that continuously identifies the correct model parameters μ for UP, when given a short recent history of the states and actions.

The training process can be formulated as a supervised learning problem with the input being a history rollout H and the output being the model parameters μ under which the input rollout is generated:

$$\theta^* = \operatorname{argmin}_{\theta} \sum_{(H_i, \mu_i) \subseteq B} \|\phi_{\theta}(H_i) - \mu_i\|^2 \quad (1)$$

where θ are the parameters of the neural net ϕ_{θ} .

Although the training data can be entirely obtained from simulation, the amount of data can be intractably large to thoroughly cover the input space. Our key observation is that OSI only needs to be accurate for the trajectories that are likely to be observed when performing the tasks of interest. As such, we randomly sample the space of μ where UP is trained for. For each sampled $\bar{\mu}_i$, we simulate N rollouts using the policy $\pi(\mathbf{x}, \bar{\mu}_i)$ and under the dynamics $f_{\bar{\mu}_i}$. We then generate short history segments from each rollout and store them in the training buffer B (Line 3-13, Algorithm 2).

After optimizing ϕ using Equation 1, we found that the performance of the combined system, UP-OSI, was much worse than simply using UP given true model parameters $\bar{\mu}$. This result is not surprising (retrospectively) because our OSI has only “seen” the motion sequences generated by a control policy $\pi(\mathbf{x}, \bar{\mu})$ under a dynamic model $f_{\bar{\mu}}$ where their model parameters are consistent. In other words, all the training examples so far only cover the “good cases” where the control policy is operating optimally under a given dynamic model. When we tested OSI with an unseen initial sequence, OSI was likely to make some error in the prediction. This error is exacerbated because the next sequence that OSI will see is generated by a control using an erroneously predicted $\hat{\mu}$ under the true model parameters $f_{\bar{\mu}}$, where $\hat{\mu} \neq \bar{\mu}$.

Algorithm 1 Learning UP

```

1: Randomly initialize UP network  $\pi$ 
2: for  $i = 1 : K$  do
3:   Initialize rollout buffer  $R$ 
4:    $\mu \sim \rho_{\mu}$ 
5:    $\mathbf{x} \sim \rho_0$ 
6:   while  $R.size \leq MaxStep$  do
7:      $\mathbf{u} = \pi(\mathbf{x}, \mu)$ 
8:      $\mathbf{x} = f_{\mu}(\mathbf{x}, \mathbf{u})$ 
9:      $r, terminated = Reward(\mathbf{x}, \mathbf{u})$ 
10:    Push  $(\mathbf{x}, \mathbf{u}, r)$  into  $R$ 
11:    if  $terminated$  then
12:       $\mu \sim \rho_{\mu}$ 
13:       $\mathbf{x} \sim \rho_0$ 
14:  Update  $\pi$  with data in  $R$  using TRPO
return  $\pi$ 

```

Our solution is to iteratively improve OSI by introducing “bad cases” with mismatched model parameters used for control ($\pi(\mathbf{x}, \hat{\mu})$) and for forward simulation ($f_{\hat{\mu}}$). For each iteration, we generate more training examples using the current OSI and UP. We randomly sample in the space of μ and generate rollouts like before. However, we feed the μ predicted by the current OSI into UP, instead of the true model parameters $\bar{\mu}$ used for forward simulation. Note that in Line 25 of Algorithm 2, the dynamic model has the parameter $\bar{\mu}$ which is different from the one used for the control (Line 22, 24). Mixing the mismatched training examples with previously generated ones, we train OSI again using Equation 1. After a small number of iterations (3-5, see Section IV), the performance of UP-OSI becomes close to the performance of UP that is provided with the true model parameters.

IV. EVALUATION

We evaluate UP-OSI on four motor control problems. In each example, the control policy does not know the true model parameters in advance and relies on OSI to identify the parameters during execution. We vary different model parameters, such as mass, friction coefficient, or task-related parameters to demonstrate that UP-OSI can successfully perform all the motor skills under unknown dynamic models. We compare the performance of UP-OSI against the performance of the condition, *UP-true*, which uses UP given the true model parameters. The performance of UP-true can be regarded as an informal upper bound for UP-OSI.

All results presented in this work¹ are simulated in PyDart2 [2], a python wrapper for DART [1], which is a multi-body physics simulator supported by Gazebo. The simulation timestep is set to 0.002s. For UP, we use a neural network with two hidden layers, comprised of 64 units in both hidden layers with tanh activation function, followed by a linear fully connected final layer. For OSI, we use three hidden layers, with 256, 128, and 64 hidden units and tanh activation function.

¹Supplementary video at <https://youtu.be/MWf6mcwPSW0>

Algorithm 2 Learning OSI

```
1: Randomly initialize OSI network  $\phi$ 
2: Initialize training buffer  $B$ 
3: for  $i = 1 : K$  do
4:    $\bar{\mu} \sim \rho_\mu$ 
5:   for  $j = 1 : N$  do
6:     Initialize history queue  $H$ 
7:     Fill  $H$  by simulating under  $\pi(\mathbf{x}, \bar{\mu})$  and  $f_{\bar{\mu}}$ 
8:     for  $t = 0 : T - 1$  do
9:       Pop  $H$ 
10:       $\mathbf{u}_t = \pi(\mathbf{x}_t, \bar{\mu})$ 
11:       $\mathbf{x}_{t+1} = f_{\bar{\mu}}(\mathbf{x}_t, \mathbf{u}_t)$ 
12:      Push  $(\mathbf{x}_{t+1}, \mathbf{u}_t)$  in  $H$ 
13:      Store  $(H, \bar{\mu})$  in  $B$ 
14: Optimize  $\phi$  using data in  $B$ 
15: while not converge do
16:   for  $i = 1 : K$  do
17:      $\bar{\mu} \sim \rho_\mu$ 
18:     for  $j = 1 : N$  do
19:       Initialize history queue  $H$ 
20:       Fill  $H$  by simulating under  $\pi(\mathbf{x}, \bar{\mu})$  and  $f_{\bar{\mu}}$ 
21:       for  $t = 0 : T - 1$  do
22:          $\hat{\mu} = \phi(H)$ 
23:         Pop  $H$ 
24:          $\mathbf{u}_t = \pi(\mathbf{x}_t, \hat{\mu})$ 
25:          $\mathbf{x}_{t+1} = f_{\bar{\mu}}(\mathbf{x}_t, \mathbf{u}_t)$ 
26:         Push  $(\mathbf{x}_{t+1}, \mathbf{u}_t)$  in  $H$ 
27:         Store  $(H, \bar{\mu})$  in  $B$ 
28:   Optimize  $\phi$  using data in  $B$ 
return  $\phi$ 
```

We add a dropout layer for OSI after each hidden layer with a dropout rate of 0.1.

The learning process for UP takes 500 iterations of TRPO updates. The amount of data collected during each iteration varies by the difficulty of the tasks. We run five iterations for training OSI. At each iteration we sample 30 different μ values and collect 5 seconds for each μ . The motion history length is a hyperparameter that can be tailored for different tasks. It represents a trade-off between the prediction accuracy of the model parameters and the speed that the policy can adapt to changing environments. We found that $H = 3$ works well for all our experiments.

A. Double inverted pendulum with unknown center of mass

We begin with a classic motor control problem: balancing a double inverted pendulum. We define the reward function as,

$$r(\mathbf{x}) = -k_1(\sigma_1 + \sigma_2)^2 - k_2|p_{cart}| + 10,$$

where σ_1 and σ_2 are angles of the two poles from the upright configuration, p_{cart} is the position of the cart and k_1 and k_2 are the corresponding weights of the two terms. We normalize the angles to be in $[0, \pi]$ and use $k_1 = 10.0$, $k_2 = 1.0$ in our experiment. The length of the two poles are both $0.5m$. We terminate the simulation when $|p_{cart}| \geq 5$ or $(\sigma_1 + \sigma_2) \geq 0.5\pi$.

The unknown model parameter for this problem is the center of mass of the lower pole, which has an unknown offset, $(\mu, 0.2\mu)$, from the geometric center. To ensure that the control policy would need to apply different strategies to balance the pendulum when different model parameters are given, we allow the offset to vary across a wide range: $\mu \in [-0.6m, 0.6m]$. Note that the purpose of the vertical offset (0.2μ) is to break the symmetry of the problem to further increase the difficulty of control.

At each training iteration of UP, we collect 150,000 samples using the physics simulator. Figure 2(b) shows the normalized performance of the trained UP-OSI across different μ values, comparing against UP-true (the informal upper bound). The performance of each μ in Figure 2(b) is the normalized average accumulated reward of 20 rollouts starting from a randomly perturbed initial state. If the performance value is above 1.0, the double inverted pendulum is able to balance. We also compare UP-OSI to a policy with conventional state input and control output but trained by data simulated from a range of model parameters (denoted as “regular” in Figure 2(b)). The purpose of this comparison is to show that providing the model parameters as input to UP results in a more powerful control policy under a range of dynamic models.

To demonstrate the learning process of OSI network, we plot the same reward-model parameter graph at each iteration of training. As shown in Figure 2(c), the performance of UP-OSI improves over time and approaches the performance of UP-true. In Figure 2(d), we plot the mean and the standard deviation for the model parameter identified by the trained OSI for each ground truth μ on x-axis. This shows that OSI is indeed able to identify the model parameter in this task.

B. Manipulator with unknown object mass

In this example, we train a robot arm to grab a block and throw it to a certain height but not beyond. The arm is initially pointing down and the block is in the air near the gripper of the arm. Similar motor skill can be observed in the serving of a tennis ball. The state \mathbf{x} includes the joint position \mathbf{q} , joint velocity $\dot{\mathbf{q}}$ of the robot arm, and the position of the block \mathbf{p}_{block} . The reward function is defined as:

$$r(\mathbf{x}, \mathbf{u}) = -k_1 r_h - k_2 \|\mathbf{u}\|^2 - k_3 \|\dot{\mathbf{q}}\|^2 + 35$$
$$r_h = \begin{cases} h_{target} - h_{block}, & \text{if } h_{block} \leq h_{target} \\ 0, & \text{otherwise} \end{cases},$$

where $k_1 = 10$, $k_2 = 1e - 5$, $k_3 = 1e - 3$, $h_{target} = 2m$ and h_{block} is the height of the block. We terminate the rollout when the box falls below $-0.2m$ or when the block is more than $0.8m$ away horizontally. By giving zero reward beyond h_{target} , we encourage the robot arm to throw the block in a way that it has low velocity when it reaches h_{target} , such that it can stay in the high reward region as long as possible. The unknown model parameters is the mass of the block. The robot needs to infer the weight of the block and use the right amount of effort to throw it up to the right height.

During the training of UP, we collect 50,000 samples for each iteration. The performance of UP-OSI and UP-true is

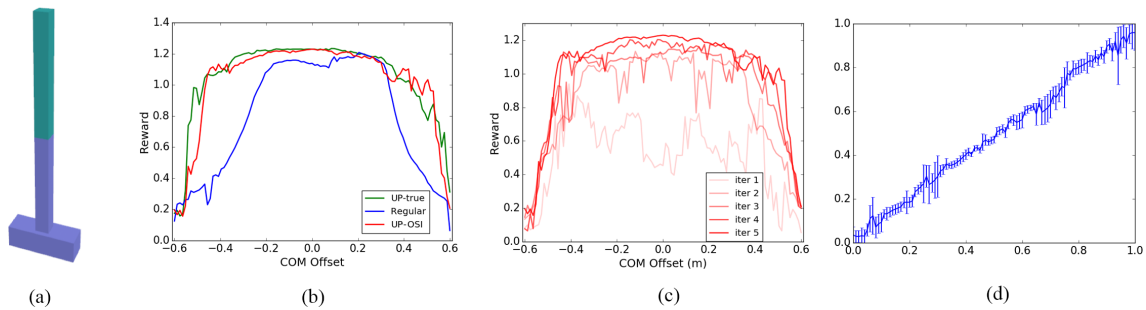


Fig. 2. Results for the Double Inverted Pendulum Task. (a) Illustration of the task. (b) Performance of UP-true, "Regular" Controller and UP-OSI. The horizontal axis is the model parameter (center of mass), and the vertical axis is the performance. (c) The evolution of optimizing UP-OSI. The reward across the training range of μ increases iteratively. (d) Mean and standard deviation of the predicted model parameter. The x-axis indicates the true model parameters while the y-axis indicates the predicted ones by OSI. The model parameters have been normalized to be in $[-1, 1]$

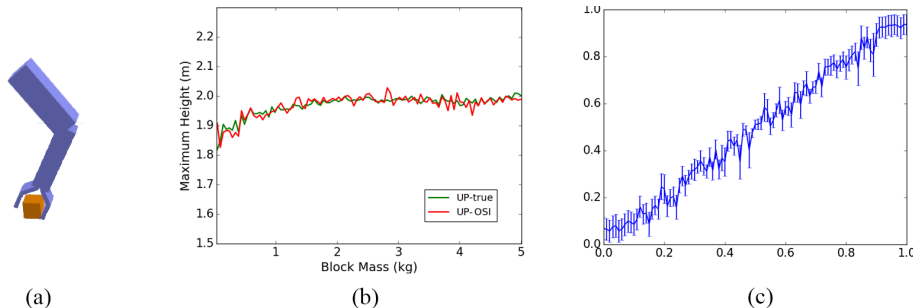


Fig. 3. Results for the Robot Arm object throwing task. (a) Illustration of the Task. (b) Performance of UP-true and UP-OSI. The horizontal axis is the model parameter (mass of the object), and the vertical axis is the performance (maximum height of the object). (c) Mean and standard deviation of the predicted model parameter. The x-axis indicates the true model parameters while the y-axis indicates the predicted ones by OSI. The model parameters have been normalized to be in $[-1, 1]$.

plotted in Figure3(b). We measure the performance by the highest point reached by the block. The closer to $h_{target} = 2m$, the better the performance. We also plot the mean and standard deviation of the predicted block mass throughout the test, as shown in Figure3(c).

C. Hopper with unknown friction coefficient

Correctly identifying contact information is crucial to many locomotion tasks. In this example, we demonstrate that our method can be applied to identify the friction coefficient at the contact point in an online fashion. The task is to control a single leg robot in 2D, the Hopper, to hop forward as fast as possible without falling. The reward is defined as

$$r(\mathbf{x}, \mathbf{u}) = k_1 \dot{\mathbf{x}} - k_2 \|\mathbf{u}\|^2 + 3.0,$$

where $k_1 = 1, k_2 = 0.002$ for our experiments. The unknown model parameters is the friction coefficient with the range $\mu \in [0.3, 1.0]$. We plot the maximum distance traveled by the hopper before the termination criteria is satisfied (the hopper falls or the maximum length of the rollout is reached) instead of the reward value to better visualize the performance of the hopper. The input to UP include the joint position of the hopper \mathbf{q} , the joint velocity $\dot{\mathbf{q}}$ and the friction coefficient μ between the foot and the ground. Note that we don't use position in the forward direction in the input state, because it is not directly related to the task.

We use 75,000 samples each iteration during the training of UP. Figure 4(b) shows the performance of UP-OSI compared

to UP-true. Due to the difficulty of the task, UP can only perform well around $\mu = [0.6, 1.0]$. However, being able to identify friction coefficients in this range, i.e. between the coefficient for wood-concrete contact and that for rubber-concrete contact, is sufficient for most practical applications. Figure 4(c) shows the mean and standard deviation of the predicted model parameter during the test at each μ .

D. Cart-pole swing-up with unknown pole length and unknown attached mass

To solve the classic cart-pole swing-up problem, the control policy needs to learn not only how to balance the pole, but also how to swing it up from a straight down position. Our experiment makes two modifications to increase the difficulty of the problem. First, we limit the force used by the cart to be within $[-40N, 40N]$. As such, the controller must swing the pole back and forth before it rises up. We also attach an additional mass to the tip of the pole to mimic the weight lifting task (Figure 5(a)).

We use a variant of the reward function suggested by [15]: $r_\sigma = w\sigma^2 + v \log(\sigma^2 + a)$, where σ is the angle of the pole. The first term encourages fast learning of swing-up motion and the second term encourages fast learning of balance. In our experiment, we set $w = 1, v = 1, a = 0.1$. Similar to the double inverted pendulum task, we also add a term to encourage the cart to stay at the center of the track, $r_{cart} = |p_{cart}|$. Together, our reward function is defined as:

$$r(\mathbf{x}) = -k_1 r_\sigma - k_2 r_{cart} + 10.0,$$

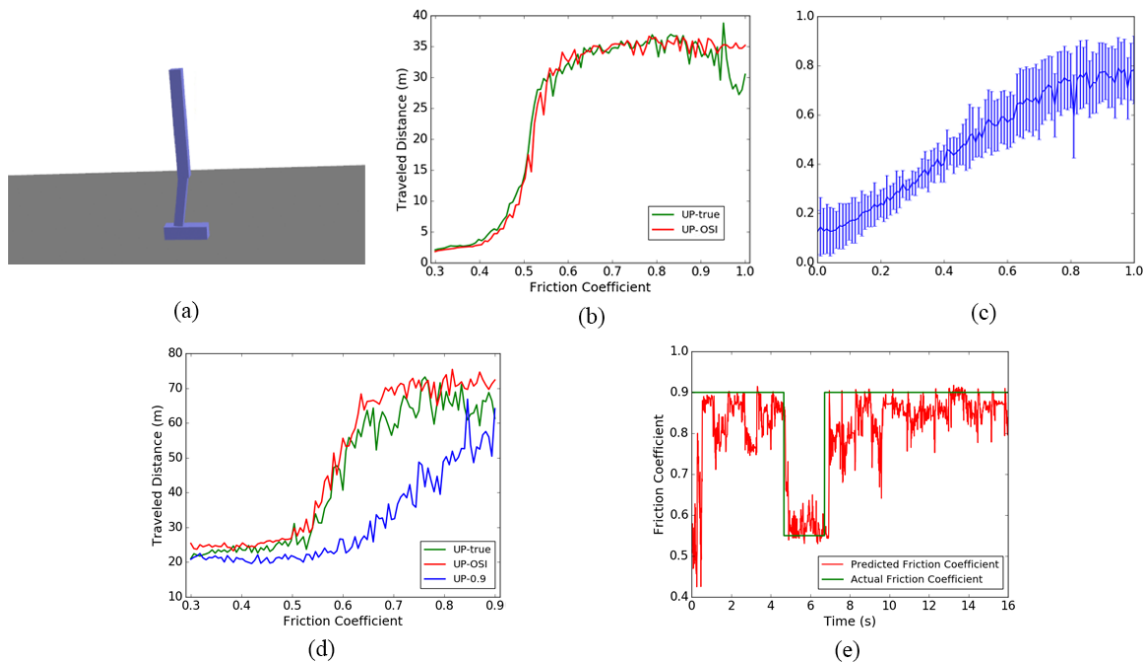


Fig. 4. Results for the Hopper task. (a) Illustration of the task. (b) Performance of UP-true and UP-OSI. The horizontal axis is the model parameter (friction coefficient), and the vertical axis is the performance (maximum distance traveled in 1,000 simulation steps). (c) Mean and standard deviation of the predicted model parameter. The x-axis indicates the true model parameters while the y-axis indicates the predicted ones by OSI. The model parameters have been normalized to be in $[-1, 1]$. (d) Performance of varying contact friction test. We tested UP-true, UP-OSI and UP with input friction coefficient fixed at 0.9 for 2,000 simulation steps. (e) OSI-predicted and actual friction coefficient in varying contact friction test.

where $k_1 = 1.0, k_2 = 0.2$. In our setup, the pole has 0 position when it is upright. During the simulation, we randomly initialize the position of the pole to be either π or $-\pi$ with a small noise drawn from $\mathcal{N}(0, 0.005)$. We terminate the rollout when the pole rotates more than 4π from the initial position, or when the cart is more than $2m$ from the center.

The unknown model parameters in this example includes the additional mass attached to the top of the pole ($\mu_{mass} \in [0.1kg, 1.0kg]$), and the length of the pole ($\mu_{length} \in [0.2m, 0.8m]$). To closely compare with the cart-pole examples in [12], where they control an inverted pendulum with varying pole length using a RNN with the whole history trajectory as input, we train OSI to estimate the velocity of the system, instead of directly giving the true velocity to the policy as part of the state. As such, the space of model parameters for this example is \mathbb{R}^4 . At each iteration of UP training, we run 70,000 samples. We normalize the resulting reward such that if a policy achieves averaged accumulated reward more than one, then it usually can swing up and balance the cart-pole system. Figure 5(b) and (c) show that UP-OSI can achieve high reward for a range of unknown pole lengths, similar to the inverted pendulum result shown in [12], but UP-OSI only requires three time steps of history as input. In addition, the mass attached to the tip is also an unknown that needs to be simultaneously identified with the pole length.

E. Generalization to varying model parameter

We run the trained UP-OSI for Hopper on a track with varying friction coefficients to test its generalizability. We create a track with friction coefficient $\mu_{const} = 0.9$ everywhere except

for the region between $20m$ to $30m$. We then vary the friction coefficient μ_{vary} in this region and plot the performance of the controller with regard to μ_{vary} . Figure 4(d) shows the performance of UP-OSI and the UP-true. Note that UP-true was given the ground truth friction coefficient at each time instance as if it has a perfect contact-friction sensor on the foot, while UP-OSI needs to identify this information based on the recent history of the motion. The results show that UP-OSI can achieve comparable and sometimes better performance than UP-true. We also test the performance of UP with fixed input $\mu = 0.9$, i.e. the hopper isn't aware of the change in friction coefficient, which is shown as the blue curve in Figure 4(d). The worse performance shows that it is crucial to detect the varying friction coefficient in order to success in this task.

In Figure 4(e), we plot the friction coefficient predicted by OSI over time for one specific $\mu_{vary} = 0.55$. We can see that OSI can identify the changes in model parameter during the task. Note that we did not provide any training examples with temporally-varying μ when training either UP-true or UP-OSI networks.

F. Generalization beyond training range

Another aspect of the policy generalizability is whether it works with model parameters that were not seen during the training phase. We perform such test on the cart-pole swing up problem with a pole length range of $[0.8m, 1.4m]$, which is 100% beyond the original training range. We also linearly increase the the attached mass with the pole length. We test both UP-true and UP-OSI with this extended range, which is shown in Figure 5(d). The result shows that UP-OSI can work

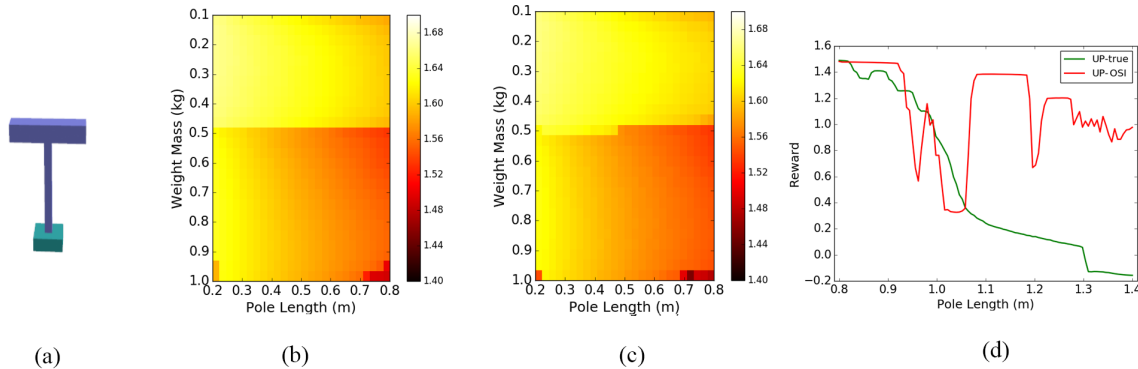


Fig. 5. Results on the Cart-Pole Swing-Up task. (a) Illustration of the task. (b) Performance of UP-true visualized in the domain of pole length and weight mass. (c) Performance of UP-OSI. (d) Performance with model parameters exceeding the training range 100%.

for a large range of unknown pole length and attached mass with only position information as input. More interestingly, UP-OSI significantly outperforms UP-true in this range of μ unseen during training.

V. DISCUSSION

While UP-OSI has demonstrated a wide range of success, we recognize a few limitations that require further investigation. First, identifying high-dimensional model parameters remains an unsolved challenge. Although we demonstrated that OSI can perform well for model parameters in \mathbb{R}^4 (the cart-pole example), more rigorous analysis is required to evaluate the sample-efficiency of UP-OSI for high-dimensional model parameter space. In all our experiments, we train OSI with five iterations and the resulting UP-OSI can achieve similar performance as the baseline. However, a theoretical upper bound and convergence conditions are not established in this work. Further, our current implementation assumes both the policy and the dynamic model are deterministic, but UP-OSI can easily be extended to a stochastic formulation.

This work primarily focuses on identifying model parameters unrelated to uncertainty in sensors and actuators partly because many of these issues have been addressed in the literature of robust control. One important source of uncertainty in a dynamic model is latency. The nature of latency might complicate the scheme of history queue used in OSI, but it is nevertheless an important and challenging problem to tackle.

An alternative approach to UP-OSI is to train an “end-to-end” control policy that takes as input a sequence of motion history and directly outputs the optimal control, rather than explicitly decoupling the process into system identification and control. Indeed, previous work [12] has shown that a recurrent network can learn to control a dynamic system with unknown model parameters. We chose to design a control policy that explicitly takes in the model parameters as input, because this additional information in the input allows the policy to specialize for each model parameter, which can potentially improve the performance. In addition, we conjecture that decoupling a large network trained by reinforcement learning into two moderate-sized networks, one of which is trained

by supervised learning, might reduce the learning time and improve the sample-efficiency.

For near future, we would like to train UP-OSI using one simulator but test it on a different one. If successful, our next step will be to test UP-OSI on real-world robotic platforms. One potential issue of our method is that the real-world dynamics cannot be well represented by the chosen parameterization of the dynamic models. However, UP-OSI allows the dynamic model to be time-varying, resulting in a more expressive parametric model to match the real-world dynamics. Essentially, the algorithm modifies its estimate of the dynamic model at every time instance based on recent real-world observations. This rapid adaptation of dynamic model allows the control policy to make prompt correction as it is also time-varying.

VI. CONCLUSION

We have introduced a new approach to creating motion policies that are formed by coupling a Universal Policy together with an On-line System Identification model. A key aspect of such policies is that they are robust under a wide variety of dynamic models, and indeed, they determine the dynamic model parameters on-the-fly. This approach was created with the goal of developing control policies for real-world tasks by performing extensive off-line training using physical simulation and deep reinforcement learning. While we have yet to transfer such policies to real robots, their performance gives several promising indications. First, the UP-OSI control policies give almost the same performance as the UP-true baseline, while having to deduce the unknown dynamic models on their own. Second, UP-OSI can perform *better* than the Universal Policy alone, in the case where the dynamic model is outside of the training range. Finally, by constantly estimating the dynamic model at every time-step, such control policies can adapt to a changing environment such as locomotion over a surface with varying friction.

ACKNOWLEDGMENTS

We thank Charles C. Kemp, Alexander Clegg, Zackory Erickson and Ariel Kapusta for their help with this work. This work is supported by NSF award IIS-1514258.

REFERENCES

- [1] DART: Dynamic Animation and Robotics Toolkit". URL <http://dartsim.github.io/>.
- [2] Pydart2. URL <https://github.com/sehoonha/pydart2>.
- [3] Pieter Abbeel and Andrew Y. Ng. Exploration and Apprenticeship Learning in Reinforcement Learning. In *International Conference on Machine Learning*, pages 1–8, 2005.
- [4] Pieter Abbeel, Morgan Quigley, and Andrew Y. Ng. Using Inaccurate Models in Reinforcement Learning. In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, pages 1–8, 2006.
- [5] Josh C. Bongard and Hod Lipson. Nonlinear System Identification Using Coevolution of Models and Tests. *IEEE Trans. Evolutionary Computation*, 9(4):361–384, 2005. URL <http://dx.doi.org/10.1109/TEVC.2005.850293>.
- [6] Paul Christiano, Zain Shah, Igor Mordatch, Jonas Schneider, Trevor Blackwell, Joshua Tobin, Pieter Abbeel, and Wojciech Zaremba. Transfer from Simulation to Real World through Learning Deep Inverse Dynamics Model. *arXiv preprint arXiv:1610.03518*, 2016.
- [7] Bruno Da Silva, George Konidaris, and Andrew Barto. Learning parameterized skills. *arXiv preprint arXiv:1206.6398*, 2012.
- [8] Marc Deisenroth and Carl E Rasmussen. PILCO: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pages 465–472, 2011.
- [9] Michel Gevers. System identification without lennart ljung: what would have been different? 2006.
- [10] Shixiang Gu, Timothy Lillicrap, Zoubin Ghahramani, Richard E Turner, and Sergey Levine. Q-Prop: Sample-Efficient Policy Gradient with An Off-Policy Critic. *arXiv preprint arXiv:1611.02247*, 2016.
- [11] Sehoon Ha and Katsu Yamane. Reducing Hardware Experiments for Model Learning and Policy Optimization. *IROS*, 2015.
- [12] Nicolas Heess, Jonathan J Hunt, Timothy P Lillicrap, and David Silver. Memory-based control with recurrent neural networks. *arXiv preprint arXiv:1512.04455*, 2015.
- [13] Stephen James and Edward Johns. 3D Simulation for Robot Arm Control with Deep Q-Learning. *arXiv preprint arXiv:1609.03759*, 2016.
- [14] Sylvain Koos, Jean-Baptiste Mouret, and Stéphane Doncieux. Crossing the reality gap in evolutionary robotics by promoting transferable controllers. In *Genetic and Evolutionary Computation Conference*. ACM, 2010. ISBN 978-1-4503-0072-8. URL <http://doi.acm.org/10.1145/1830483.1830505>.
- [15] Sergey Levine, Nolan Wagener, and Pieter Abbeel. Learning contact-rich manipulation skills with guided policy search. In *Robotics and Automation (ICRA), 2015 IEEE International Conference on*, pages 156–163. IEEE, 2015.
- [16] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-End Training of Deep Visuomotor Policies. *Journal of Machine Learning Research*, 17(39): 1–40, 2016.
- [17] Sergey Levine, Peter Pastor, Alex Krizhevsky, and Deirdre Quillen. Learning Hand-Eye Coordination for Robotic Grasping with Deep Learning and Large-Scale Data Collection. *CoRR*, abs/1603.02199, 2016. URL <http://arxiv.org/abs/1603.02199>.
- [18] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [19] 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):529–533, 2015.
- [20] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy P Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning*, 2016.
- [21] Igor Mordatch, Kendall Lowrey, and Emanuel Todorov. Ensemble-CIO: Full-body dynamic motion planning that transfers to physical humanoids. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2015.
- [22] S. Nolfi and D. Floreano. *Evolutionary Robotics: The Biology, Intelligence, and Technology*. MIT Press (Cambridge, MA), 2000.
- [23] Xue Bin Peng, Glen Berseth, and Michiel van de Panne. Terrain-Adaptive Locomotion Skills Using Deep Reinforcement Learning. *ACM Transactions on Graphics*, 35(4), 2016.
- [24] Karime Pereida, Rikky R. P. R. Duivenvoorden, and Angela P. Schoellig. High-precision trajectory tracking in changing environments through L1 adaptive feedback and iterative learning. In *Proc. of the IEEE International Conference on Robotics and Automation (ICRA)*, 2017.
- [25] Lerrel Pinto and Abhinav Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 3406–3413. IEEE, 2016.
- [26] Ali Punjani and Pieter Abbeel. Deep learning helicopter dynamics models. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 3223–3230. IEEE, 2015.
- [27] Aravind Rajeswaran, Sarvjeet Ghotra, Sergey Levine, and Balaraman Ravindran. EPOpt: Learning Robust Neural Network Policies Using Model Ensembles. *arXiv preprint arXiv:1610.01283*, 2016.
- [28] Stephane Ross and J. Andrew Bagnell. Agnostic system identification for model-based reinforcement learning. In *ICML*, 2012.

- [29] Andrei A Rusu, Matej Vecerik, Thomas Rothörl, Nicolas Heess, Razvan Pascanu, and Raia Hadsell. Sim-to-real robot learning from pixels with progressive nets. *arXiv preprint arXiv:1610.04286*, 2016.
- [30] John Schulman, Sergey Levine, Philipp Moritz, Michael I Jordan, and Pieter Abbeel. Trust region policy optimization. *CoRR, abs/1502.05477*, 2015.
- [31] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- [32] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin A. Riedmiller. Deterministic Policy Gradient Algorithms. In *ICML*, 2014.
- [33] Freek Stulp, Gennaro Raiola, Antoine Hoarau, Serena Ivaldi, and Olivier Sigaud. Learning compact parameterized skills with a single regression. In *Humanoid Robots (Humanoids), 2013 13th IEEE-RAS International Conference on*, pages 417–422. IEEE, 2013.
- [34] István Szita, Bálint Takács, and András Lörincz. ε -MDPs: Learning in varying environments. *Journal of Machine Learning Research*, 3(Aug):145–174, 2002.
- [35] Juan Cristobal Zagal, Javier Ruiz-del-Solar, and Paul Vallejos. Back-to-Reality: Crossing the Reality Gap in Evolutionary Robotics. In *IAV 2004: Proceedings 5th IFAC Symposium on Intelligent Autonomous Vehicles*. Elsevier Science Publishers B.V., 2004. ISBN 008-044237-4.