

Learning from the Hindsight Plan – Episodic MPC Improvement

Aviv Tamar¹, Garrett Thomas¹, Tianhao Zhang¹, Sergey Levine¹, Pieter Abbeel^{1,2,3}

Abstract—Model predictive control (MPC) is a popular control method that has proved effective for robotics, among other fields. MPC performs re-planning at every time step. Re-planning is done with a limited horizon per computational and real-time constraints and often also for robustness to potential model errors. However, the limited horizon leads to suboptimal performance. In this work, we consider the iterative learning setting, where the same task can be repeated several times, and propose a policy improvement scheme for MPC. The main idea is that between executions we can, offline, run MPC with a longer horizon, resulting in a *hindsight plan*. To bring the next real-world execution closer to the hindsight plan, our approach learns to re-shape the original cost function with the goal of satisfying the following property: short horizon planning (as realistic during real executions) with respect to the shaped cost should result in mimicking the hindsight plan. This effectively consolidates long-term reasoning into the short-horizon planning. We empirically evaluate our approach in contact-rich manipulation tasks both in simulated and real environments, such as peg insertion by a real PR2 robot.

I. INTRODUCTION

Model predictive control (MPC), also known as receding horizon control, is an effective model-based control method that is well-suited for complex dynamical systems, with a wide range of applications in robotics (see, e.g., [1], [2]) among other domains. In MPC, a model of system dynamics, either learned or known, is used to plan a locally optimal control policy for a *limited horizon*, starting from the current state of the system. The first control in the plan is executed, and then re-planning is performed from the new state.

When the system dynamics are not known exactly, as is often the case in practice, MPC can be integrated with online system identification to simultaneously learn dynamics and plan controls. This approach has been successfully applied in various robotic tasks, from aerial vehicle flight [3], [4], to contact-rich object manipulation [5], [6]. In these applications, the limited horizon of MPC serves a dual purpose: maintaining tractability of the planning problem and mitigating error propagation during planning as a result of inaccurate models [7]. Indeed, MPC has been successfully applied to problems with challenging dynamics, such as cutting vegetables [5] and object manipulation [6].

One drawback of MPC is that planning with a limited horizon can lead to suboptimal policies. For example, consider a task of navigating an environment with an obstacle. The MPC policy can only maneuver around the obstacle when it is within the planning horizon. If the robot were to encounter the same task several times, we should naturally expect it to *learn* to maneuver around the obstacle even before it enters the planning horizon. However, state-of-the-art MPC methods that only learn dynamics, such as [5] and

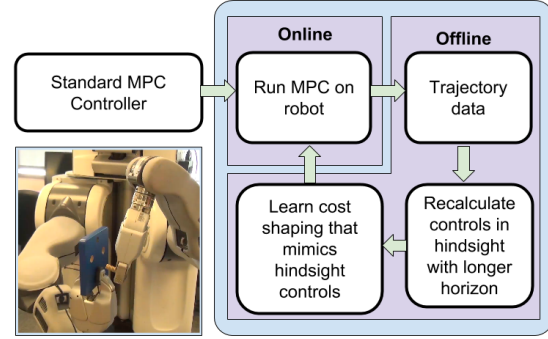


Fig. 1: Right: A schematic of our method. After running an MPC controller with a limited horizon, we perform an offline calculation of the MPC on the same trajectory but with a longer horizon. This *hindsight plan* is used to learn a cost shaping for the next execution of the MPC. Left: PR2 learning to insert a wooden peg into a hole.

[6], are prone to repeatedly produce suboptimal behavior in each episode of the task due to the limited horizon.

In this work, we propose a method that improves the MPC policies in episodic tasks. Our main insight is that, between episodes, we can revisit the MPC planning computations that were performed online, and recompute them *offline* with a longer horizon, and with potentially better dynamics. The difference between this *hindsight plan* and the actions that were actually performed can be used to drive a policy improvement *between episodes*. In particular, we learn a neural network cost shaping for MPC using supervised learning, by minimizing this difference. The result is a method that incorporates long-term reasoning into MPC, while maintaining the benefits of short-horizon planning.

In comparison to previous policy improvement methods, which are typically based on value functions or policy gradients [8], our approach exploits the predictive nature of MPC to drive policy improvement, by contrasting the predicted actions with actions calculated in hindsight. This allows us to circumvent the difficulties of value function approximation and the sample inefficiency of policy gradients.

We show that our method effectively improves MPC policies for contact-rich manipulation tasks, such as peg insertion, in both simulated and real environments.

II. RELATED WORK

A standard approach for mitigating the limited planning horizon in MPC is to introduce a value function as the terminal cost, also known as *infinite horizon MPC* [9], [10]. When the model is not fully known, or too large to calculate the value function, approximation techniques must be employed. However, it has been observed that standard value function approximation methods perform poorly when combined with MPC [11]. Zhong et al. [11] also proposed

¹EECS Department, UC Berkeley, ²ICSI, UC Berkeley, ³OpenAI

an improved approximation based on state discretization. However, that method does not scale to high-dimensional systems such as the 7 DoF robot arm considered in our experiments. We note that the cost shaping in our method can be set to have the form of a terminal value function. However, instead of learning directly from observed scalar costs, as typical for value function approximation [8], our method learns the cost shaping from the hindsight plan, which contains *desired controls* and hence is much more informative.

Our episodic formulation of improving the MPC controller is similar to the setting of iterative learning control (ILC) [12], [13], [14], [15]. However, in ILC, the goal is typically to follow a reference trajectory, which is not suitable for the manipulation tasks we consider here. Our formulation allows the goal to be specified as a general optimal control problem. We note that a very recent work [16] proposed an extension of ILC that does not require a reference trajectory, and also uses a form of learning MPC, based on adding a value function to the MPC cost, which is calculated for all previously visited states. That approach, however, requires on a controller that can deterministically drive the system to any previously visited state, which does not hold for the real-world manipulation tasks we consider here.

Reinforcement planning [17] considers a discrete high-level planner for a continuous low-level control algorithm, and learns cost parameters of the planner using reinforcement learning methods. Our approach, in comparison, exploits the fact that our controller is based on MPC to learn from the hindsight plan, and is not limited to discrete planning. We also note that our method is different from *hindsight optimization* [18], in which simulation is used to calculate a value function online.

Learning a cost function from observed controls is a form of inverse optimal control (IOC) [19], [20]. While in IOC, the actions are assumed to be generated by an expert demonstrator, our method does not require such demonstrations, and the actions used for learning the cost shaping are generated by a planning algorithm which learns the model from interaction, as in model-based reinforcement learning (RL) [21], [22] and adaptive control [23]. Model-based RL is a standard approach in robotic control [8], achieving state-of-the-art results in various domains, from helicopter flight [24] to contact-rich manipulation [25]. The use of MPC in RL allows for extremely data-efficient learning [6], which is especially important in robotics, where robot interaction time is often costly. In this work we improve the MPC controllers of [6] for contact-rich manipulation tasks. However, our method is general, and can be applied to any adaptive MPC algorithm, and any domain for which MPC is suitable.

III. PRELIMINARIES

In this work, we consider the standard episodic reinforcement learning or optimal control setting in discrete time with episode length T . We denote by \mathbf{x}_t the state at time t , and \mathbf{u}_t the control. The goal is to generate a control sequence $\mathbf{U}_{0:T} = (\mathbf{u}_0, \dots, \mathbf{u}_T)$ that minimizes the total cost $\sum_{t=0}^T \ell_t(\mathbf{x}_t, \mathbf{u}_t)$ given an initial state \mathbf{x}_0 , where

ℓ_t is a predefined task-specific scalar cost function, under dynamics given by $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t)$. At time step t , adaptive MPC methods calculate the control \mathbf{u}_t by first estimating the dynamics model and then approximately solving for the optimal action.

Concretely, let $\mathcal{D}_t = \{(\mathbf{x}_s, \mathbf{u}_s, \mathbf{x}_{s+1})\}_{s=0}^{t-1}$ denote a dataset of the observed interaction with the system up to time t . Let $H \leq T$ denote the MPC horizon, which is a parameter of the algorithm. At time t , MPC¹ first uses \mathcal{D}_t to predict system dynamics \hat{f}_s^t for the next H time steps, i.e. for $t \leq s \leq t+H$. Using the predicted dynamics, MPC then calculates the H -horizon optimal actions $\mathbf{U}_{t:t+H}^* = (\mathbf{u}_t^*, \dots, \mathbf{u}_{t+H}^*)$, by solving

$$\begin{aligned} \arg \min_{\mathbf{u}_t, \dots, \mathbf{u}_{t+H}} \quad & \sum_{s=t}^{t+H} \ell_s(\mathbf{x}_s, \mathbf{u}_s), \\ \text{s.t.} \quad & \mathbf{x}_{s+1} = \hat{f}_s^t(\mathbf{x}_s, \mathbf{u}_s), \quad \forall s = t, \dots, t+H. \end{aligned} \quad (1)$$

The action $\mathbf{u}_t = \mathbf{u}_t^*$ is then taken, and the system transitions to a new state $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t)$, from which the MPC optimization is repeated. When the episode ends, at $t = T$, the system is reset to the initial condition \mathbf{x}_0 and the process iterates. To reduce clutter, we omit the episode index from the notation.

IV. THE HINDSIGHT PLAN

In principle, if we know the true system dynamics f , and set H equal to the episode length T , the solution of Eq. (1) would lead to an optimal policy. However, in practice, both the computational burden of solving Eq. (1), and the error in dynamics prediction, necessitate the use of a shorter horizon $H < T$ (and often much shorter $H \ll T$) for online MPC optimization, which results in a *suboptimal* policy. To mitigate this sub-optimality, we propose to revisit the planning computation *after an episode has ended*, using a longer horizon, and potentially better dynamics, as we describe in this section. This additional planning computation, which we term the *hindsight plan*, will allow us to improve our MPC policy, as we shall later describe.

Let \bar{H} denote the horizon of the hindsight plan, where typically² $\bar{H} \geq H$. The *hindsight action* $\bar{\mathbf{u}}_t$ at time t is defined as $\bar{\mathbf{u}}_t = \bar{\mathbf{u}}_t^*$, where $\bar{\mathbf{U}}_{t:t+\bar{H}}^* = (\bar{\mathbf{u}}_t^*, \dots, \bar{\mathbf{u}}_{t+\bar{H}}^*)$ is the solution to the following *hindsight planning problem*:

$$\begin{aligned} \arg \min_{\bar{\mathbf{u}}_t, \dots, \bar{\mathbf{u}}_{t+\bar{H}}} \quad & \sum_{s=t}^{t+\bar{H}} \ell_s(\mathbf{x}_s, \bar{\mathbf{u}}_s), \\ \text{s.t.} \quad & \mathbf{x}_{s+1} = \hat{f}_s^s(\mathbf{x}_s, \bar{\mathbf{u}}_s), \quad \forall s = t, \dots, t+\bar{H}, \end{aligned} \quad (2)$$

where the initial state \mathbf{x}_t is the state observed at time t during the original MPC execution, and the terms \hat{f}_s^s are the same terms that were already calculated for the original MPC planning problem (1). Henceforth, we collectively term \mathbf{x}_t and \hat{f}_s^s as the MPC trajectory information $I \doteq \{\mathbf{x}_t, \hat{f}_s^s\}$.

Let us emphasize the differences between the hindsight planning problem (2) and the original MPC planning problem (1). First, we plan with a longer horizon \bar{H} . Second, at each time step t , we use the dynamics predictions $\hat{f}_t^t, \dots, \hat{f}_{t+\bar{H}}^{t+\bar{H}}$

¹Note that this slightly non-standard notation of the MPC algorithm makes explicit the online dynamics learning.

²With perfect dynamics, setting $\bar{H} = T$ would be optimal. However, with an inaccurate dynamics model, setting $\bar{H} < T$ may be preferred.

that were calculated at times $t, \dots, t + \bar{H}$ of the online MPC execution, as opposed to the predictions that were calculated at time t , $\hat{f}_t^t, \dots, \hat{f}_{t+\bar{H}}^t$ in Eq. (1). These predictions were not available at time t during the online MPC execution, as they use observations *from later time steps* in the episode³.

The main assumption underlying our approach is that the hindsight plan produces improved actions compared to the original MPC execution. This assumption is motivated by the improved dynamics prediction in (2), which is based on future data and, more importantly, the much longer planning horizon. Thus, in the hindsight plan, the two main sources of MPC sub-optimality are removed. We also note that calculating the hindsight plan typically requires more computation than the original MPC planning during the episode execution. However, the real-time requirement of MPC can be relaxed, because it can be computed offline, and the calculation can be performed concurrently and in parallel for all t .

V. MPC POLICY IMPROVEMENT

In this section, we show how the hindsight plan can be used to drive a policy improvement algorithm, by using it for learning a *cost shaping* for MPC. The idea is to learn a cost shaping that encourages the online MPC solution to be *more similar to the solution in hindsight*.

Let us revisit the MPC optimization problem, and add to the original cost a cost shaping term, $\delta_s(\mathbf{x}_s, \mathbf{u}_s, \theta)$, parametrized by some vector θ ,

$$\begin{aligned} \arg \min_{\mathbf{u}_t, \dots, \mathbf{u}_{t+H}} \sum_{s=t}^{t+H} \ell_s(\mathbf{x}_s, \mathbf{u}_s) + \delta_s(\mathbf{x}_s, \mathbf{u}_s, \theta), \\ \text{s.t. } \mathbf{x}_{s+1} = \hat{f}_s^t(\mathbf{x}_s, \mathbf{u}_s), \quad \forall s = t, \dots, t+H \end{aligned} \quad (3)$$

where, similarly to Eq. (2), the trajectory information $I = \{\mathbf{x}_t, \hat{f}_s^t\}$ contains the terms that were calculated for the original MPC planning problem (1). The cost-shaping term $\delta_s(\mathbf{x}_s, \mathbf{u}_s, \theta)$ can be represented, for example, by a neural network parametrized by θ . We denote by $\mathbf{u}_t(\theta)$ the first action in the solution of (3) for time t . That is, $\mathbf{u}_t(\theta)$ denotes the action MPC *would have taken* during the episode, had the cost-shaping parameter been θ . We also denote by $\bar{\mathbf{u}}_t^0$ the first action in the solution of (3) with $\delta_s(\mathbf{x}_s, \mathbf{u}_s, \theta) \equiv 0$, that is, the solution with no cost shaping applied.

We aim to learn a parameter θ that encourages the MPC solution of Eq. (3) to be similar to the the hindsight plan (2). We therefore propose to learn θ by minimizing the similarity loss \mathcal{L} , defined as

$$\mathcal{L}(\theta) = \sum_{t=0}^T \|\mathbf{u}_t(\theta) - \bar{\mathbf{u}}_t^0\|^2 + \lambda \|\mathbf{u}_t(\theta) - \mathbf{u}_t^0\|^2, \quad (4)$$

where $\lambda \geq 0$ acts as a regularization term that controls the change from the online MPC policy. Note that the dependence of $\mathbf{u}_t(\theta)$ on θ can be quite complex, as it encapsulates the solution of the planning problem (3). A similar problem is encountered in IOC [19], [20], and indeed, our approach can be seen as performing IOC with the hindsight plan replacing

the expert demonstration. For planning algorithms that can be represented as a computation graph [26], such as linear quadratic regulator (LQR) and value iteration [27], the loss \mathcal{L} can be minimized efficiently using gradient based methods, such as L-BFGS [28]. We also note that our method can be used to learn additional parameters of the planning algorithm, such as the dynamics and discount factor, by simply adding them to the parameter vector θ . In this work, however, we focus on learning cost shaping.

Furthermore, if $\delta_s(\mathbf{x}_s, \mathbf{u}_s, \theta)$ is set to 0 for all $s < t + H$, the cost shaping becomes a terminal value function, which is exactly the infinite horizon MPC formulation [9], [10]. A crucial difference in our approach, however, is that we learn the parameters θ by minimizing the similarity loss \mathcal{L} , which directly measures the difference between the online and hindsight actions. Standard value function methods first approximate a value function, and then approximately solve the planning problem with that value function, resulting in two sources of approximation error, possibly explaining their poor performance observed in previous work [11].

The solution to Eq. (4) provides us with the parameters of a shaping cost that causes the MPC controller to mimic the hindsight controller, and results in a single step of policy improvement. However, once we run the MPC with the shaped cost in the system, we obtain a new trajectory, which we can perform hindsight planning on. This results in an iterative policy improvement algorithm, which we term *hindsight iterative MPC* (HIMPC). Let $i \geq 0$ denote an iteration of the algorithm, and let θ_i denote the cost-shaping parameters used for the MPC execution at iteration i , where we assume that the first iteration $i = 0$ is executed without any cost shaping. We denote by $\mathcal{L}_i(\theta)$ the similarity loss optimization problem (4), where the trajectory information I_i is calculated from execution of the MPC controller at iteration i , with cost-shaping parameter θ_i . We calculate θ_{i+1} by minimizing $\sum_{k=0}^i \mathcal{L}_k(\theta)$, which assures that we aggregate the hindsight information obtained at previous iterations.

The general HIMPC algorithm is summarized in Algorithm 1. So far, we have not discussed the specific algorithms used for dynamics predictions and planning in MPC (1), nor the functional form of the cost shaping, and in principle, HIMPC can be combined with any adaptive MPC method such as [5] and [6]. However, the specific choice of algorithm and cost shaping will determine the computational complexity of solving the cost-shaped MPC planning (3) in real time. In addition, the tractability of minimizing the similarity loss (4) depends on the specific algorithm and cost shaping structure. In the next section, we describe a particular implementation of HIMPC based on LQR planning and Gaussian Mixture Model (GMM) dynamics. We also propose a neural-network (NN) cost shaping form that is both efficient to optimize and interpretable.

VI. AN LQR IMPLEMENTATION OF HIMPC

In this section we describe a particular implementation of the HIMPC algorithm, based on LQR planning and GMM dynamics learning. We are inspired by the work of [6], in which similar methods were used within MPC for effectively performing contact-rich manipulation tasks. In addition, we

³Potentially, we can use the data set at the end of the episode \mathcal{D}_T to re-estimate the dynamics at every time step, and obtain even better dynamics predictions. However, we found it sufficient to use the dynamics that were already calculated during the online MPC execution.

Algorithm 1: HIMPC Algorithm

- 1 Run original MPC controller (1), and collect trajectory information I_0
 - 2 Calculate hindsight plan (2) with I_0
 - 3 Find cost-shaping that mimics hindsight controller
 $\theta_0 = \arg \min_{\theta} \mathcal{L}_0(\theta)$
 - 4 **for** $i=1, 2, \dots$ **do**
 - 5 Run MPC controller with cost-shaping parameter
 θ_{i-1} , and collect trajectory information I_i
 - 6 Calculate hindsight plan (2) with I_i
 - 7 Solve $\theta_i = \arg \min_{\theta} \sum_{k=0}^i \mathcal{L}_k(\theta)$
 - 8 **end**
-

propose a novel cost shaping structure, that is both efficient to optimize, and has an intuitive interpretation as learning ‘way points’.

The online dynamics adaptation algorithm of [6] is based on a Bayesian approach for estimation, where the recent observations within an episode are used to estimate a linear dynamics model, using observations from previous episodes as a Bayesian prior. This produces, for every time step t , a time-varying linear dynamics model of the form $\mathbf{x}_{s+1} = \hat{f}_s^t(\mathbf{x}_s, \mathbf{u}_s) = A_s^t \mathbf{x}_s + B_s^t \mathbf{u}_s$, for $s = t, \dots, t+H$. The method is further described in Appendix I.

As the loss function ℓ_s , we use a quadratic loss of the form $\ell_s(\mathbf{x}_s, \mathbf{u}_s) = (\mathbf{x}_s - \mathbf{x}^*)^\top Q_s (\mathbf{x}_s - \mathbf{x}^*) + \mathbf{u}_s^\top R_s \mathbf{u}_s$, where \mathbf{x}^* is some goal state of the system, and Q_s and R_s are a positive semi-definite and positive definite matrices, respectively⁴. Such a cost function is standard for many control tasks [29], and is particularly suitable for the manipulation experiments we consider, where the task is specified as moving the robot end-effector to some goal position, such as pushing a peg into a hole.

With linear dynamics and a quadratic loss, the MPC planning problem (1) becomes a standard LQR problem [29], for which a solution can be calculated efficiently by dynamic programming, as described Appendix II. The mapping from the trajectory information $\{\mathbf{x}_t, \hat{f}_s^t\}$ to the action \mathbf{u}_t in this case can be written as the sequence of matrix multiplications and inversions in the LQR solution. This mapping can be represented as a computation graph, and the gradient $\partial \mathbf{u}_t / \partial \mathbf{x}^*$ can be easily calculated using modern automatic differentiation packages such as Theano [30].

We are now ready to introduce our cost shaping formulation. At time t , for $s = t, \dots, t+H$, we consider a cost shaping of the form:

$$\delta_s(\mathbf{x}_s, \mathbf{u}_s, \theta) = g(\mathbf{x}_t, \theta)^\top Q_s \mathbf{x}_s, \quad (5)$$

where $g(\mathbf{x}_t, \theta)$ is a neural network that has the current state \mathbf{x}_t as its input, θ as its weights, and a vector-valued output with the same dimensionality as \mathbf{x}_t . Note that adding a linear term to a quadratic form is equivalent to changing the center of the quadratic, up to a constant. Therefore, the shaped cost can be written as:

$$\begin{aligned} \ell_s(\mathbf{x}_s, \mathbf{u}_s) + \delta_s(\mathbf{x}_s, \mathbf{u}_s, \theta) = \\ (\mathbf{x}_s - \hat{\mathbf{x}}^*(\mathbf{x}_t, \theta))^\top Q_s (\mathbf{x}_s - \hat{\mathbf{x}}^*(\mathbf{x}_t, \theta)) + \mathbf{u}_s^\top R_s \mathbf{u}_s + \text{const}, \end{aligned} \quad (6)$$

⁴Extending our approach to non-quadratic loss functions is straightforward, by using a second-order Taylor expansion. See, e.g., [2] for details.

where $\hat{\mathbf{x}}^*(\mathbf{x}_t, \theta) = \mathbf{x}^* + g(\mathbf{x}_t, \theta)$ is a *modified goal position*. Thus, our cost shaping has the intuitive interpretation of modifying the goal position, which can be thought of as learning state-dependent way points. Note that g depends on the current state observation \mathbf{x}_t , and not on \mathbf{x}_s , therefore solving LQR with the cost in (6) is equivalent to solving the original LQR, just with a different \mathbf{x}^* , making the solution similarly tractable.

In addition, calculating the gradient $\partial \mathbf{u}_t / \partial \theta = \partial \mathbf{u}_t / \partial \hat{\mathbf{x}}^* \cdot \partial \hat{\mathbf{x}}^* / \partial \theta$ is also tractable. This gradient can be used for minimizing the similarity loss (4) with standard optimization algorithms such as L-BFGS [28]. Note that the gradient $\partial g / \partial \mathbf{x}$ is not required since g depends on \mathbf{x}_t , which is part of the trajectory information I .

An Illustrative Example

In this section we discuss an application of HIMPC to a simple 2D navigation task with obstacles. The goal of this example is to illustrate the function of the hindsight plan, and the learned cost shaping. Further quantitative analysis of this experiment is presented in Section VIII.

Consider the 2D navigation task depicted in Figure 2. A particle with mass needs to navigate to a goal position by using vertical and horizontal forces. The particle may collide with the impenetrable gray colored obstacles, which when in contact, apply normal and frictional forces to the particle.

Ideally, when starting from the initial position shown in Figure 2, the particle should navigate to the opening between the obstacles and from there continue to the goal. Such a plan would minimize the total cost in this domain. However, when a MPC policy is applied, the controller first navigates to the obstacle location which is closest to the goal, as shown in the dashed-red line in Figure 2, since its limited horizon planning (here $H = 10$) and imperfect dynamics model fail to take into account the future obstacle collision. Only after experiencing the obstacle, the controller navigates alongside it towards the goal. While the initial MPC policy is able to solve the task, its solution is clearly sub-optimal.

The wide red and black arrows in Figure 2 show the action that MPC selected, \mathbf{u}_t , and the action that the hindsight plan (with $\bar{H} = 30$ in this case) has chosen, $\bar{\mathbf{u}}_t$, at a particular time step. Due to the better dynamics prediction and longer horizon, the hindsight plan correctly predicts the future collision, and takes an appropriate action.

The solid-black line in Figure 2 shows the trajectory performed by the HIMPC algorithm, after 5 episodes of learning. The controller learned to first navigate to the opening, as desired. In addition, we visualize the cost shaping by plotting a black quiver plot of the direction to the modified goal position $\hat{\mathbf{x}}^*(\mathbf{x}, \theta) - \mathbf{x}$. Note how the modified goal position at the beginning of the trajectory orients towards the opening, and note the difference with the red quiver plot, which shows the direction to the original goal $\mathbf{x}^* - \mathbf{x}$.

VII. PRACTICAL IMPROVEMENTS OF HIMPC

As with most machine learning algorithms, and neural networks in particular, a successful application of the method requires some technical know-how [31], [32]. In this section

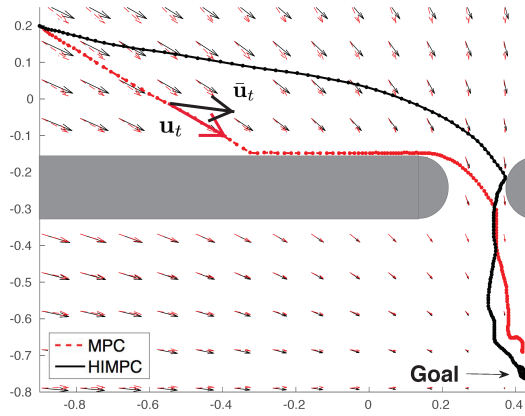


Fig. 2: 2D navigation between obstacles. The task is to navigate a particle in 2D between the gray obstacles to the goal position. Dashed-red line: MPC controller. Solid-black line: HIMPC controller. Wide arrows show the action MPC chose at a particular time \mathbf{u}_t , and the corresponding hindsight action $\bar{\mathbf{u}}_t$. The quiver plot shows the direction to the (shaped) goal position. Note that for HIMPC, the modified goal position for states at the top left is shifted to the right, leading to a trajectory that correctly navigates into the opening between the obstacles.

we report several technical procedures that we found to improve the performance of HIMPC in our experiments.

a) Add control noise to the MPC: This helps the MPC controller get around ‘local minima’ in the trajectory by random exploration, and the HIMPC then learns a cost shaping that consolidates this trajectory improvement. In particular, we used the exploration scheme of [6] in our experiments.

b) Collect several trajectories in each iteration: We found that collecting several roll-outs of the same shaped MPC controller at each iteration (lines 1 and 5 in Algorithm 1) stabilizes the neural network training.

c) Wait for successful MPC runs before initiating cost-shaping: When the standard MPC fails in the task due to a bad initial dynamics model, its trajectories do not contain enough useful knowledge for the hindsight planning. We therefore wait until several successful trajectories occur before starting HIMPC. We found measuring success by a threshold on the final distance to the target to perform well. While this condition is task dependent, in our experiments we did not find it to be sensitive to the threshold magnitude.

d) Early stopping of cost-shaping The learned cost shaping can potentially direct the controller to a different goal position than \mathbf{x}^* . This often happens in early iterations, when not enough successful trajectories have been observed, and can cause the algorithm to destabilize, as it no longer receives successful trajectories. A solution to this problem is to turn off the cost shaping (during the run) if the shaped cost has converged, while the original cost has not. Such a fix guarantees that HIMPC is as stable as the original MPC.

VIII. EXPERIMENTS

In this section, we experimentally evaluate the HIMPC algorithm in simulated and real robot experiments. In this work, we focus on contact-rich manipulation, following the work of Fu et al. [6], though our method can be applied to

other tasks where MPC is applicable. In particular, we focus on various insertion tasks, which are important for assembly, and for which the improvement of HIMPC over standard MPC can be easily visualized.

In our evaluation, we aim to answer the following two questions:

- 1) Can HIMPC improve upon standard MPC?
- 2) Can HIMPC improve upon a standard episodic model-based RL approach (i.e., without MPC)?

For our simulations, in addition to the 2D navigation task discussed in Section VI, we consider two variants of a peg insertion task. All our simulations were performed using the MuJoCo physics engine [33], and our code, which is based on the guided policy search repository [34], will be made available. Additionally we experiment on a peg insertion task with a real PR2 robot. In both simulated and real experiments, we apply direct torque control at 20Hz. For the simulated experiments, the state space is 26-dimensional, consisting of the positions and angular velocities of the 7 joints, and positions and velocities of 2 points on the end-effector. For the real robot experiment, an additional end-effector point was added, resulting in 32 state dimensions.

In our evaluation, we compare HIMPC to the MPC method of [6], and, as a baseline, to iLQG – a state-of-the-art model-based RL method [22]. Our dynamics model uses a GMM prior, combined with online dynamics adaptation [6], as further described in Appendix I. The iLQG method of [22] uses a similar GMM for a dynamics prior, but combines with time-varying linear dynamics, and therefore constitutes a fair comparison. The same quadratic cost function was used for all algorithms. We note that the optimization in iLQG is performed over the full episode, and should in principal converge to a (locally) optimal solution. However, MPC is expected to be much more sample-efficient, as was demonstrated in [6].

In terms of computation time, hindsight planning demands were comparable to system reset times between episodes, and had negligible effect. However, solving the optimization problem in (4) was much slower, requiring several minutes of computation, since our Theano-based implementation [30] cannot differentiate a matrix inverse for multiple samples in parallel. We expect future automatic differentiation packages to substantially reduce the computation time.

A. Simulated 2D Navigation

This task, as depicted in Figure 2, requires moving a particle through an opening towards a goal position in 2D space. The available controls apply horizontal and vertical forces to the particle, and obstacles effect friction and normal forces. Episodes are 200 time-steps long, and three trajectories are executed at each iteration. The horizon H and hindsight horizon \bar{H} were chosen as 10 and 30, respectively. The cost-shaping neural network g consists of two fully connected hidden layers, each with 25 units and tanh activations. The inputs to g are the particle position and velocity.

The dynamics GMM prior was initialized from a single episode with a random initial policy. In addition, after each episode we update the prior with the samples from the most

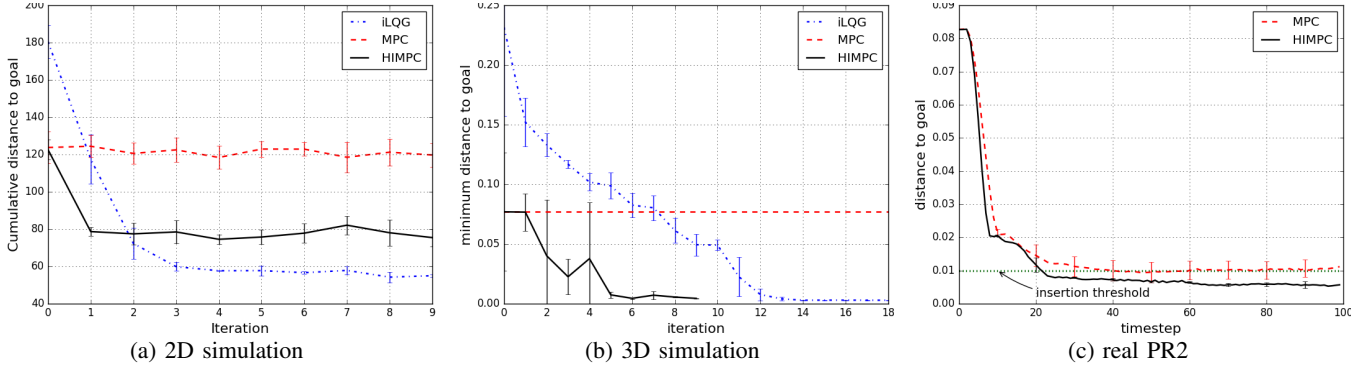


Fig. 3: Experimental results. (a): Performance plots for the 2D navigation task. (b): Performance plots for the simulated 3D peg insertion task. We show the minimal distance from the target, averaged over the trajectories in the episode. HIMPC significantly improves the baseline MPC performance, and converges faster than iLQG. (c): Performance plots for the peg insertion task with the real PR2. HIMPC inserts the peg faster, and with better precision.

recent episode, to potentially improve the baseline MPC algorithm between iterations.

Our performance evaluation is the cumulative distance to the goal over the entire episode, which measures the speed of reaching the goal. In Figure 3a, we plot the averaged performance over four runs using different random seeds. As may be observed, HIMPC significantly improves over standard MPC. Interestingly, updating the dynamics prior did not significantly improve the original MPC (as can be seen by relatively similar performance across episodes). This is since in this simple domain, the dynamics prior of the first episode is already good enough. Note that the baseline iLQG reaches a better solution than MPC. This is expected, as iLQG solves the problem with the full horizon and therefore requires more samples than MPC.

Generalization and the effect of \bar{H} . The HIMPC algorithm learns to improve performance when iteratively starting the task from the same initial position. An important consideration, however, is its sensitivity to the change of initial position. In this experiment we demonstrate that the neural network cost shaping can generalize to initial positions that were not trained on.

Using the simulated 2D domain described in section VIII-A, we ran HIMPC using samples collected from five different initial positions. We then evaluated performance when starting from 16 different test positions, arranged in a 4×4 grid. We evaluated HIMPC for various values of \bar{H} , as well as standard MPC, on the test positions, and present the results (averaged across 10 test trials) in Figure 5. Observe that as \bar{H} increases, the performance of HIMPC improves consistently for all test positions. Furthermore, when $\bar{H} > H$, HIMPC outperforms standard MPC, even when run from initial positions that were not seen during training. This result shows that the benefit of planning with a longer horizon is indeed captured by the learned cost shaping, and also demonstrates the generalization capability of the neural network.

B. Simulated 3D Peg Insertion

In this task, a 7-Dof robot arm must insert a cylindrical peg into a tight-fitting rectangular hole, as depicted in Figure 4a & 4b. The controls are the torques applied to the 7 joints, and

the state space consists of the positions and angular velocities of these joints and the 3D positions and velocities of 2 points on the end-effector, which are hereafter referred to as the *EE points*. The goal is specified by the coordinates of the EE points (but not the joints), when the peg is fully inserted in the hole. Thus, there is a rotational degree of freedom in the goal position, since there are only 2 EE points. This non-trivial task requires solving both a kinematic problem, and a control problem with complex contact dynamics, and has been used as a benchmark in previous studies [35].

Episodes are 400 time-steps long, and 3 trajectories are executed at each iteration. The horizon H is 10, while the hindsight horizon \bar{H} was chosen as 60. The NN g had 2 fully connected hidden layers of sizes [100, 25] with tanh activations, and its inputs were the positions of the EE points. The dynamics GMM prior was initialized from a single episode with a random policy. In this experiment we did not update the prior after the first episode, as we observed it to decrease performance. We attribute this to the GMM method, which is sensitive to the choice of clusters and the input distribution, and can degrade performance when additional samples are added. Using dynamics models such as neural networks [6] or Gaussian processes [36] could potentially improve the dynamics learning, as we plan to investigate in future work. We emphasize, however, that improving MPC dynamics should also improve the performance of HIMPC.

To make a fair comparison with iLQG, which learns time-dependent linear dynamics models, we executed 6 trajectories of 200 time-steps at each iLQG iteration. This results in a better performance of iLQG, while the number of samples (total time steps) at each iteration is the same as HIMPC.

Our performance evaluation is the minimal distance to the goal over the episode, which measures the success of inserting the peg. In Figure 3b, we plot our results, averaged over the trajectories in the episode. As may be observed, HIMPC significantly improves over the standard MPC, which is equivalent to the performance of HIMPC in the first episode. In this domain, HIMPC solves the task with significantly fewer samples than iLQG.

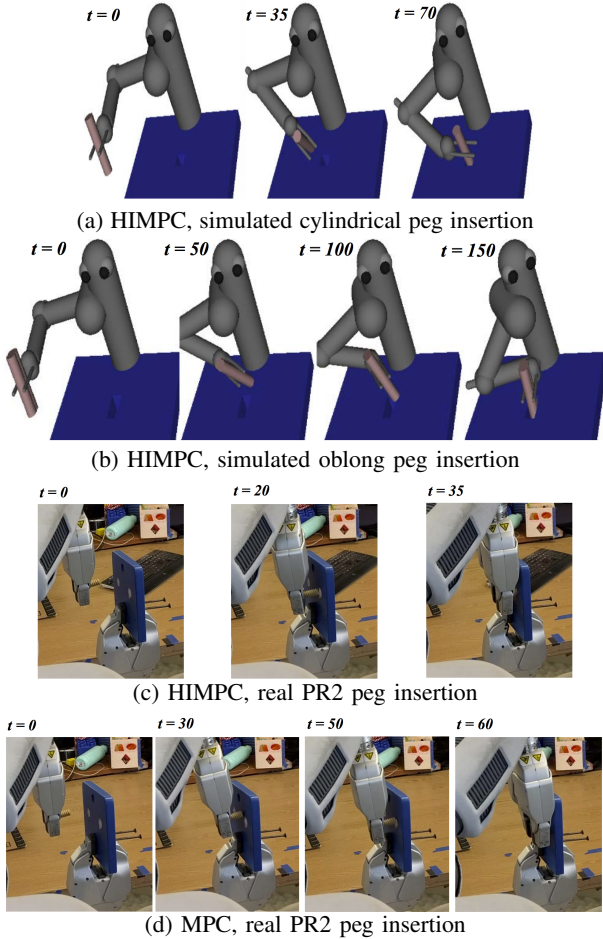


Fig. 4: Visualization of HIMPC and MPC in peg insertion tasks. (a,b): The HIMPC policy for the simulated peg insertion tasks. Note the different arm orientation for the cylindrical (a) and oblong (b) peg. (c,d) The HIMPC policy and baseline MPC policy for the real peg insertion task. Note that MPC first reaches an intermediate point on the plate, while HIMPC directly reaches the hole.

C. Simulated 3D Oblong Peg Insertion

In this experiment, we demonstrate how HIMPC with random exploration noise in the control, can learn additional structure in the task, and represent it in the shaped cost.

This task is similar to the previous peg insertion task, with a difference that the peg and hole have an oblong shape, requiring a particular orientation of the peg for a successful insertion. As before, the goal is specified by the coordinates of the 2 EE points, when the peg is fully inserted in the hole. Thus, the cost function *does not contain information* about the correct orientation for solving the task. With sufficient exploration noise, standard MPC can sometimes solve this task. Our goal is to show that HIMPC can consolidate the information from these ‘lucky’ runs, and learn a cost shaping that guides the peg into the correct orientation.

As before, episodes are 400 time-steps long, and 3 trajectories are executed at each iteration. The horizon H is 10, while the hindsight horizon \bar{H} was chosen as 60, and the NN g had 2 fully connected hidden layers of sizes [100, 25] with tanh activations. In this case, however, we added to the NN

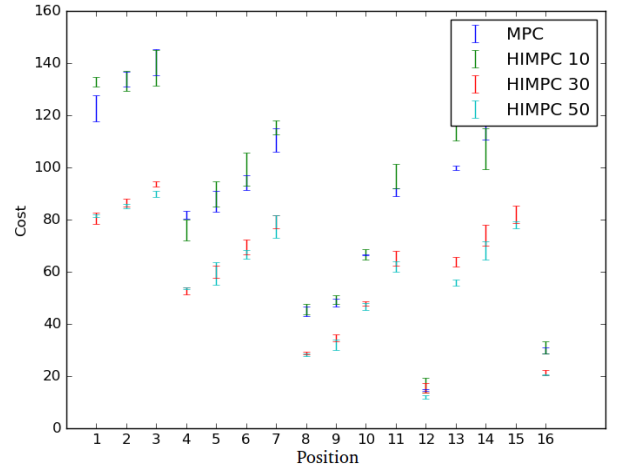


Fig. 5: Generalization and the effect of \bar{H} . Performance (cumulative distance to goal, averaged over 10 trials with standard deviation error bars) of MPC and HIMPC on the 2D navigation task, starting from 16 different initial positions, for various settings of \bar{H} . Note that performance consistently improves with a longer horizon for all test positions. These test positions were not used during training, and demonstrate the generalization capability of the neural network.

shaping-cost an additional EE point, which was not used for the original cost specification. Thus, the cost-shaping has the capacity to represent an orientation of the peg. We emphasize that this additional EE point *was not* part of the MPC cost, and the orientation can only be learned from the hindsight plan on ‘lucky’ trajectories.

In Table I we report the success rate of peg insertion for MPC, and HIMPC after 6 episodes of learning. We also report on success or failure when no control noise is added. After 6 episodes, HIMPC has learned a cost shaping that orients the peg correctly, and succeeds in the task, even without control noise. MPC on the other hand, cannot solve the task without exploration noise. This behavior can be further visualized in the supplemental video⁵.

TABLE I: Success Rates for Oblong Peg Insertion

	W/ Noise	W/O Noise
MPC	4/10	Fail
HIMPC	8/10	Success

D. Real PR2 Experiments

We evaluated our method on a peg insertion task with the PR2 robot. The robot is tasked with inserting a small wooden peg into a hole in a wooden plate. The task specification is similar to the simulated peg insertion of Section VIII-B, and the goal position is specified by the position of 3 EE points, when the peg is fully inserted in the hole.

We constructed a GMM dynamics prior from 40,000 samples of actions in free space, collected by running the iLQG algorithm for 30 minutes, reaching random goal positions. Such a large data set was required for a stable MPC control. In addition, we used a shorter hindsight horizon than in simulation $\bar{H} = 30$, to account for the less accurate dynamics.

⁵<https://sites.google.com/site/himpchindsightplan/>

In Figure 3c, we plot the distance of the EE points from their goal position, for the original MPC controller, and for HIMPC after 3 episodes of learning, averaged over 5 executions of the controller. Similarly to the 2D task described above, the MPC controller in this task tries to approach the goal position in a straight line, and bumps into the wooden plate. It then glides into the opening to insert the peg. HIMPC, on the other hand, learns to directly insert the peg into the opening, and therefore reaches the goal faster. This behavior can be visualized in the supplemental video⁵.

IX. CONCLUSION

In this work we introduced a new approach for policy improvement in repeated tasks. Rather than using value functions or policy gradients – the traditional drivers of policy improvement in RL – our method employs an online MPC policy, and suggests improved actions based on an offline hindsight calculation once an episode has ended.

We demonstrated a significant improvement over standard MPC on several complex manipulation tasks with contacts, and a notable improvement in sample efficiency over state-of-the-art model based RL.

In future work we intend to investigate the use of different dynamics prediction models in our method, and applications in different robotics domains such as quadrotors. In addition, the explicit use of the prediction error as a driver for policy improvement could potentially be used in different RL algorithms.

APPENDIX I

Dynamics Prediction: We used the dynamics prediction method of [6]. An exponential moving average of the observations is maintained throughout the episode $\hat{\mu}_t \leftarrow \beta \hat{\mu}_{t-1} + (1 - \beta) \mathbf{p}_t$, where $\mathbf{p}_t = [\mathbf{x}_{t-1}; \mathbf{u}_{t-1}; \mathbf{x}_t]$ is the t^{th} observation and β is a discounting factor that causes the model to forget old data. A similar exponential moving covariance is maintained by $\hat{\Sigma}_t \leftarrow \beta \hat{\Sigma}_{t-1} + (1 - \beta) \mathbf{p}_t \mathbf{p}_t^\top$. These within-episode dynamics are mixed with a prior model of dynamics (μ_p, Σ_p) , by $\mu = \alpha_1 \hat{\mu} + (1 - \alpha_1) \mu_p$, and $\Sigma = \alpha_2 \Sigma_p + \alpha_3 \hat{\Sigma} + \alpha_4 (\mu_p - \hat{\mu})(\mu_p - \hat{\mu})^\top$, where the prior is a GMM fit to samples from previous episodes, and the mixing coefficients $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ are described in [6]. A linear dynamics model $\mathbf{x}_{t+1} = A_t \mathbf{x}_t + B_t \mathbf{u}_t$ is obtained by assuming a multivariate Gaussian distribution for $[\mathbf{x}_t; \mathbf{u}_t; \mathbf{x}_{t+1}]$ with mean μ and covariance Σ , and conditioning \mathbf{x}_{t+1} on $[\mathbf{x}_t, \mathbf{u}_t]$, producing a Gaussian distribution with mean $[A_t, B_t]$. We refer to [6] for the full details and theoretical motivation of this algorithm. To predict the dynamics for time $s = t+1$, the previous MPC policy (for time $t-1$) is used to predict the current action $\hat{\mathbf{u}}_t$, and the predicted next state is $\hat{\mathbf{x}}_{t+1} = A_t \mathbf{x}_t + B_t \hat{\mathbf{u}}_t$. The prior is queried for the dynamics of state $\hat{\mathbf{x}}_{t+1}$, and mixed with the current dynamics estimate as described above, to produce a linear dynamics model $\mathbf{x}_{s+1} = A_s^t \mathbf{x}_s + B_s^t \mathbf{u}_s$. This process can be repeated iteratively for $s = t+2, \dots, t+H$, producing a time-varying linear dynamics model for the horizon H .

APPENDIX II

LQR: Consider a linear dynamics system $\mathbf{x}_{t+1} = f_t(\mathbf{x}_t, \mathbf{u}_t) \doteq A_t \mathbf{x}_t + B_t \mathbf{u}_t$, and a quadratic loss function

$\ell_t(\mathbf{x}_t, \mathbf{u}_t) = (\mathbf{x}_t - \mathbf{x}^*)^\top D_t (\mathbf{x}_t - \mathbf{x}^*) + \mathbf{u}_t^\top R_t \mathbf{u}_t$. The Q-function and value function are both quadratic, given by

$$V(\mathbf{x}_t) = \frac{1}{2} \mathbf{x}_t^\top V_{\mathbf{x}, \mathbf{x}_t} \mathbf{x}_t + \mathbf{x}_t^\top V_{\mathbf{x}t} + \text{const}$$

$$Q(\mathbf{x}_t, \mathbf{u}_t) = \frac{1}{2} [\mathbf{x}_t; \mathbf{u}_t]^\top Q_{\mathbf{xu}, \mathbf{xu}_t} [\mathbf{x}_t; \mathbf{u}_t] + [\mathbf{x}_t; \mathbf{u}_t]^\top Q_{\mathbf{xu}t} + \text{const},$$

Using dynamics programming [29]:

$$Q_{\mathbf{xu}, \mathbf{xu}_t} = \ell_{\mathbf{xu}, \mathbf{xu}_t} + f_{\mathbf{xu}_t}^\top V_{\mathbf{x}, \mathbf{x}_{t+1}} f_{\mathbf{xu}_t}$$

$$Q_{\mathbf{xu}t} = \ell_{\mathbf{xu}t} + f_{\mathbf{xu}_t}^\top V_{\mathbf{x}_{t+1}}$$

$$V_{\mathbf{x}, \mathbf{x}_t} = Q_{\mathbf{x}, \mathbf{x}_t} - Q_{\mathbf{u}, \mathbf{x}_t}^\top Q_{\mathbf{u}, \mathbf{u}_t}^{-1} Q_{\mathbf{u}, \mathbf{x}_t}$$

$$V_{\mathbf{x}t} = Q_{\mathbf{x}t} - Q_{\mathbf{u}, \mathbf{x}_t}^\top Q_{\mathbf{u}, \mathbf{u}_t}^{-1} Q_{\mathbf{u}t}.$$

The optimal control law is linear, and given by $\mathbf{u}_t = k_t + K_t \mathbf{x}_t$, where $K_t = -Q_{\mathbf{u}, \mathbf{u}_t}^{-1} Q_{\mathbf{u}, \mathbf{x}_t}$ and $k_t = -Q_{\mathbf{u}, \mathbf{u}_t}^{-1} Q_{\mathbf{u}t}$.

ACKNOWLEDGMENT

This work was supported in part by Siemens, the DARPA SIMPLEX program, an NSF CAREER Award, and an ONR Young Investigator Award. Aviv Tamar was partially funded by the Viterbi Scholarship, Technion. Tianhao Zhang was supported by a UC Berkeley EECS Departmental Fellowship. The authors thank Ramu Chandra, Karthik Kappaganthu, and Juan L. Aparicio for fruitful discussions, and Justin Fu for useful advice, and for sharing his adaptive MPC code.

REFERENCES

- [1] E. F. Camacho and C. B. Alba, *Model predictive control*. Springer Science & Business Media, 2013.
- [2] T. Erez, K. Lowrey, Y. Tassa, V. Kumar, S. Koley, and E. Todorov, “An integrated system for real-time model predictive control of humanoid robots,” in *Humanoids*, 2013.
- [3] A. Aswani, P. Bouffard, and C. Tomlin, “Extensions of learning-based model predictive control for real-time application to a quadrotor helicopter,” in *ACC*, 2012, pp. 4661–4666.
- [4] G. Chowdhary, M. Mühlegg, J. P. How, and F. Holzapfel, “Concurrent learning adaptive model predictive control,” in *Advances in Aerospace Guidance, Navigation and Control*. Springer, 2013, pp. 29–47.
- [5] I. Lenz, R. Knepper, and A. Saxena, “Deepmpc: Learning deep latent features for model predictive control,” in *Robotics Science and Systems*, 2015.
- [6] J. Fu, S. Levine, and P. Abbeel, “One-shot learning of manipulation skills with online dynamics adaptation and neural network priors,” *IROS*, 2016.
- [7] N. Jiang, A. Kulesza, S. Singh, and R. Lewis, “The dependence of effective planning horizon on model accuracy,” in *AAMAS*. International Foundation for Autonomous Agents and Multiagent Systems, 2015, pp. 1181–1189.
- [8] J. Kober, J. A. Bagnell, and J. Peters, “Reinforcement learning in robotics: A survey,” *International Journal of Robotics Research*, 2013.
- [9] H. Chen and F. ALLGoWER, “A quasi-infinite horizon nonlinear model predictive control scheme with guaranteed stability,” *Automatica*, vol. 34, no. 10, pp. 1205–1217, 1998.
- [10] T. Erez, Y. Tassa, and E. Todorov, “Infinite-horizon model predictive control for periodic tasks with contacts,” *Robotics: Science and systems*, p. 73, 2012.
- [11] M. Zhong, M. Johnson, Y. Tassa, T. Erez, and E. Todorov, “Value function approximation and model predictive control,” in *2013 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL)*. IEEE, 2013, pp. 100–107.
- [12] K. L. Moore and J.-X. Xu, *Iterative learning control*. Taylor & Francis, 2000.
- [13] R. W. Longman, “Iterative learning control and repetitive control for engineering practice,” *International journal of control*, vol. 73, no. 10, pp. 930–954, 2000.
- [14] D. A. Bristow, M. Tharayil, and A. G. Alleyne, “A survey of iterative learning control,” *IEEE Control Systems*, 2006.
- [15] Y. Wang, F. Gao, and F. J. Doyle, “Survey on iterative learning control, repetitive control, and run-to-run control,” *Journal of Process Control*, vol. 19, no. 10, pp. 1589–1600, 2009.
- [16] U. Rosolia and F. Borrelli, “Learning model predictive control for iterative tasks,” *arXiv preprint arXiv:1609.01387*, 2016.

- [17] M. Zucker and J. A. Bagnell, "Reinforcement planning: RL for optimal planners," in *ICRA*. IEEE, 2012, pp. 1850–1855.
- [18] E. K. Chong, R. L. Givan, and H. S. Chang, "A framework for simulation-based network control via hindsight optimization," in *CDC*, 2000.
- [19] J. Rust, "Maximum likelihood estimation of discrete control processes," *SIAM Journal on Control and Optimization*, vol. 26, no. 5, pp. 1006–1024, 1988.
- [20] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *ICML*. ACM, 2004, p. 1.
- [21] M. Deisenroth and C. E. Rasmussen, "PILCO: A model-based and data-efficient approach to policy search," in *ICML*, 2011, pp. 465–472.
- [22] S. Levine and P. Abbeel, "Learning neural network policies with guided policy search under unknown dynamics," in *NIPS*, 2014.
- [23] K. J. Åström and B. Wittenmark, *Adaptive control*. Courier Corporation, 2013.
- [24] P. Abbeel, A. Coates, M. Quigley, and A. Y. Ng, "An application of reinforcement learning to aerobatic helicopter flight," *NIPS*, 2007.
- [25] S. Levine, N. Wagener, and P. Abbeel, "Learning contact-rich manipulation skills with guided policy search," in *ICRA*, 2015.
- [26] J. Schulman, N. Heess, T. Weber, and P. Abbeel, "Gradient estimation using stochastic computation graphs," in *NIPS*, 2015, pp. 3528–3536.
- [27] A. Tamar, Y. Wu, G. Thomas, S. Levine, and P. Abbeel, "Value iteration networks," *CoRR*, vol. abs/1602.02867, 2016.
- [28] D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Mathematical programming*, 1989.
- [29] B. D. Anderson and J. B. Moore, *Optimal control: linear quadratic methods*. Courier Corporation, 2007.
- [30] Theano Development Team, "Theano: A Python framework for fast computation of mathematical expressions," *arXiv e-prints*, vol. abs/1605.02688.
- [31] Y. A. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller, "Efficient backprop," in *Neural networks: Tricks of the trade*. Springer, 2012.
- [32] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Aistats*, vol. 9, 2010, pp. 249–256.
- [33] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 5026–5033.
- [34] C. Finn, M. Zhang, J. Fu, X. Tan, Z. McCarthy, E. Scharff, and S. Levine, "Guided policy search code implementation," 2016. [Online]. Available: <http://rll.berkeley.edu/gps>
- [35] M. Zhang, Z. McCarthy, C. Finn, S. Levine, and P. Abbeel, "Learning deep neural network policies with continuous memory states," in *ICRA*. IEEE, 2016, pp. 520–527.
- [36] Y. Pan, X. Yan, E. Theodorou, and B. Boots, "Adaptive probabilistic trajectory optimization via efficient approximate inference," *arXiv preprint arXiv:1608.06235*, 2016.