

---

# Value Iteration in Continuous Actions, States and Time

---

Michael Lutter<sup>1,2</sup> Shie Mannor<sup>1,3</sup> Jan Peters<sup>2</sup> Dieter Fox<sup>1,4</sup> Animesh Garg<sup>1,5</sup>

## Abstract

Classical value iteration approaches are not applicable to environments with continuous states and actions. For such environments, the states and actions are usually discretized, which leads to an exponential increase in computational complexity. In this paper, we propose continuous fitted value iteration (cFVI). This algorithm enables dynamic programming for continuous states and actions with a known dynamics model. Leveraging the continuous-time formulation, the optimal policy can be derived for non-linear control-affine dynamics. This closed-form solution enables the efficient extension of value iteration to continuous environments. We show in non-linear control experiments that the dynamic programming solution obtains the same quantitative performance as deep reinforcement learning methods in simulation but excels when transferred to the physical system. The policy obtained by cFVI is more robust to changes in the dynamics despite using only a deterministic model and without explicitly incorporating robustness in the optimization. Videos of the physical system are available at <https://sites.google.com/view/value-iteration>.

## 1. Introduction

Reinforcement learning (RL) maximizes the scalar rewards  $r$  by trying different actions  $u$  selected by the policy  $\pi$ . For a deterministic policy and deterministic dynamics, the discrete time optimization is described by

$$\pi^*(x) = \arg \max_{\pi} \sum_{i=0}^{\infty} \gamma^i r(x_i, u_i), \quad (1)$$

with the dynamics  $x_{i+1} = f(x_i, u_i)$ , the state  $x$  and discounting factor  $\gamma \in [0, 1)$  (Bellman, 1957; Puterman, 1994). The discrete time framework has been very successful to

---

<sup>1</sup>NVIDIA <sup>2</sup>Technical University of Darmstadt <sup>3</sup>Technion, Israel Institute of Technology <sup>4</sup>University of Washington <sup>5</sup>University of Toronto & Vector Institute. Correspondence to: Michael Lutter <michael@robot-learning.de>.

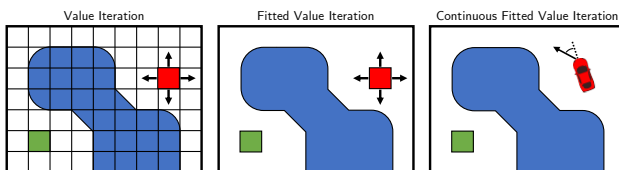


Figure 1. Value iteration (VI) can compute the optimal value function  $V^*$  and policy  $\pi^*$  for discrete state and action environments. Fitted value iteration extends VI to continuous states and discrete actions. Our proposed continuous fitted value iteration enables value iteration for continuous states and continuous actions, e.g., a car with steering rather than a car moving left, right, up and down.

learn policies for various applications including robotics (Silver et al., 2017; OpenAI et al., 2019; Da et al., 2020). However, the fixed pre-determined time discretization is a limitation when dealing with physical phenomena evolving in continuous time. For such environments, the time-step can be chosen to fit the environment. For example, in robotics the control frequency is only limited by the sampling frequencies of sensors and actuators, which is commonly much higher than the control frequency of deep RL agents. To avoid the fixed time steps, we use the continuous-time RL objective to frame an optimization that is agnostic to the time-discretization. Furthermore, this formulation simplifies the optimization as many robot dynamics are control affine at the continuous-time limit.

To solve the continuous-time RL problem, we use classical value iteration where the dynamics and reward function are known. For robotics, this setting is commonly a reasonable assumption. The reward is known as it is manually tuned to achieve the task without violating system constraints. The equations of motion and system parameters of robots are approximately known. The model is not perfect as the actuators are not ideal and the physical parameters vary slightly. However, the overall system behavior is approximately correct. The same assumption is also used within the vast sim2real literature to enable the sample efficient transfer to the physical world (Ramos et al., 2019; Chebotar et al., 2019; Xie et al., 2021).

In this paper we show for the continuous-time RL that

(1) classical value iteration (Bellman, 1957) can be extended to environments with continuous actions and states if the dynamics are control affine and the reward separable. Learning the value function  $V$  is sufficient as the policy can

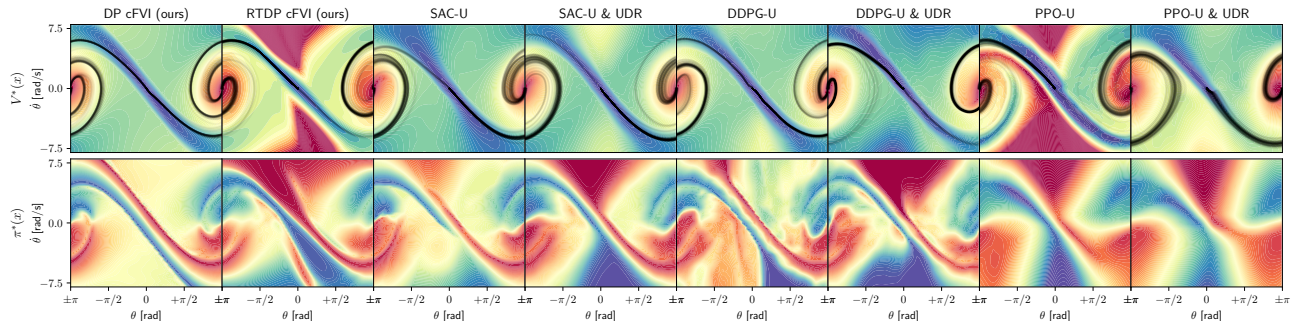


Figure 2. The optimal value function  $V^*$  and policy  $\pi^*$  and roll outs of the torque-limited pendulum computed by DP cFVI, RTDP cFVI and the deep RL baselines with uniform initial state distribution. cFVI learns a symmetric and smooth optimal policy that can swing-up the pendulum from both sides. Especially the DP variant has a very sharp ridge leading to the upward pointing pendulum. The baselines do not achieve a symmetric swing-up and usually prefer a one-sided swing-up.

be deduced from  $V$  in the continuous-time limit. Therefore, one does not need the policy optimization used by the predominant actor-critic approaches (Schulman et al., 2015; Lillicrap et al., 2015) or the discretization required by the classical algorithms (Bellman, 1957; Sutton et al., 1998).

(2) the proposed approach can be successfully applied to learn the optimal policy for real-world under-actuated systems using only the approximate dynamics model.

(3) We provide an in-depth quantitative and qualitative evaluation with comparisons to actor-critic algorithms. Using value iteration on the compact state domain obtains a more robust policy when transferred to the physical systems.

**Summary of contributions.** We extend value iteration to continuous actions as well as the extensive quantitative and qualitative evaluation on two physical systems. The latter is in contrast to prior work (Schulman et al., 2017; Haarnoja et al., 2018; Lillicrap et al., 2015; Harrison\* et al., 2017; Mandlekar\* et al., 2017), which mainly focuses on the quantitative comparison in simulation. Instead, we focus on the qualitative evaluation on the physical Furuta pendulum and Cartpole to understand the differences in performance.

In the following, we summarize continuous-time RL (Section 2) and introduce continuous fitted value iteration (Section 3). Section 4 describes suitable network representations for a Lyapunov value function. Finally, we analyze the performance on the physical systems in Section 5.

## 2. Problem Statement

We consider the deterministic, continuous-time RL problem. The infinite horizon optimization is described by

$$\pi^*(\mathbf{x}_0) = \arg \max_{\pi} \int_0^{\infty} \exp(-\rho t) r_c(\mathbf{x}_t, \mathbf{u}_t) dt \quad (2)$$

$$V^*(\mathbf{x}_0) = \max_{\mathbf{u}} \int_0^{\infty} \exp(-\rho t) r_c(\mathbf{x}_t, \mathbf{u}_t) dt \quad (3)$$

$$\text{with } \mathbf{x}(t) = \mathbf{x}_0 + \int_0^t f_c(\mathbf{x}_{\tau}, \mathbf{u}_{\tau}) d\tau \quad (4)$$

with the discounting constant  $\rho \in (0, \infty]$ , the reward  $r_c$  and the dynamics  $f_c$  (Kirk, 1970). Notably, the discrete reward and discounting can be described using the continuous-time counterparts, i.e.,  $r(\mathbf{x}, \mathbf{u}) = \Delta t r_c(\mathbf{x}, \mathbf{u})$  and  $\gamma = \exp(-\rho \Delta t)$  with the sampling interval  $\Delta t$ . The continuous-time discounting  $\rho$  is, in contrast to the discrete discounting factor  $\gamma$ , agnostic to sampling frequencies. In the continuous-time case, the Q-function does not exist (Doya, 2000).

The deterministic continuous-time dynamics model  $f_c$  is assumed to be non-linear w.r.t. the system state  $\mathbf{x}$  but affine w.r.t. the action  $\mathbf{u}$ . Such dynamics model is described by

$$\dot{\mathbf{x}} = \mathbf{a}(\mathbf{x}; \theta) + \mathbf{B}(\mathbf{x}; \theta)\mathbf{u}, \quad (5)$$

with the non-linear drift  $\mathbf{a}$ , the non-linear control matrix  $\mathbf{B}$  and the system parameters  $\theta$ . Robot dynamics models are naturally expressed in the continuous-time formulation and many are control affine. Furthermore, this special case has received ample attention in the existing literature due to its wide applicability (Doya, 2000; Kappen, 2005; Todorov, 2007). The reward is separable into a non-linear state reward  $q_c$  and the action cost  $g_c$  described by

$$r_c(\mathbf{x}, \mathbf{u}) = q_c(\mathbf{x}) - g_c(\mathbf{u}). \quad (6)$$

The action penalty  $g_c$  is non-linear, positive definite and strictly convex. This separability is common for robot control problems as rewards are composed of a state component quantifying the distance to a desired state and an action penalty. The action cost penalizes non-zero actions to avoid bang-bang control from being optimal and is convex to have a unique optimal action.

## 3. Continuous Fitted Value Iteration

First, we summarize the existing value iteration approaches and present the theory enabling us to extend value iteration to continuous action spaces. Afterwards, we introduce our proposed algorithm to learn the value function.

### 3.1. Value Iteration Preliminaries

Value iteration (VI) is an approach to compute the optimal value function for a discrete time, state and action MDP with known dynamics and reward (Bellman, 1957). This approach iteratively updates the value function of each state using the Bellman optimality principle described by

$$V^{k+1}(\mathbf{x}_t) = \max_{\mathbf{u}_{0:\ell}} \sum_{i=0}^{\ell-1} \gamma^i r(\mathbf{x}_{t+i}, \mathbf{u}_i) + \gamma^\ell V^k(\mathbf{x}_{t+\ell}),$$

with the number of look-ahead steps  $\ell$ . VI is proven to converge to the optimal value function (Puterman, 1994) as the update is a contraction described by

$$\|V_i^{k+1} - V_j^{k+1}\| \leq \gamma^\ell \|V_i - V_j\|. \quad (7)$$

For discrete actions, the greedy action is determined by evaluating all possible actions. However, VI is impractical for larger MDPs as the computational complexity grows exponentially with the number of states and actions. This is especially problematic for continuous spaces as discretization leads to an exponential increase of state-action tuples.

VI can be applied to continuous state spaces and **discrete** actions by using a function approximator instead of tabular value function. Fitted value iteration (FVI) computes the value function target using the VI update and minimizes the  $\ell_p$ -norm between the target and the approximation  $V^k(\mathbf{x}; \psi)$ . This approach is described by

$$V_{\text{tar}}(\mathbf{x}_t) = \max_{\mathbf{u}} r(\mathbf{x}_t, \mathbf{u}) + \gamma V^k(f(\mathbf{x}_t, \mathbf{u}); \psi_k) \quad (8)$$

$$\psi_{k+1} = \arg \min_{\psi} \sum_{\mathbf{x} \in \mathcal{D}} \|V_{\text{tar}}(\mathbf{x}) - V^k(\mathbf{x}; \psi)\|_p^p \quad (9)$$

with the parameters  $\psi_k$  at iteration  $k$  and the fixed dataset  $\mathcal{D}$ . The convergence proof of VI does not generalize to FVI as the fitting of the value function is not necessarily a contraction (Boyan & Moore, 1994; Baird, 1995; Tsitsiklis & Van Roy, 1996; Munos & Szepesvári, 2008). Despite these theoretical limitations, many authors proposed model-free variants using the Q-function for discrete action MDPs, e.g., fitted Q-iteration (FQI) (Ernst et al., 2005), Regularized FQI (Farahmand et al., 2009), Neural FQI (Riedmiller, 2005) and DQN (Mnih et al., 2015). The resulting algorithms were empirically successful and solved Backgammon (Tesauro, 1992) and the Atari games (Mnih et al., 2015).

For continuous actions, current RL methods use policy iteration (PI) rather than VI (Schulman et al., 2015; Lillicrap et al., 2015; ?). PI evaluates the value function of the current policy and hence, uses the action of the policy to compute the value function target. Therefore, PI circumvents the maximization required of VI (Equation 8). To update the policy, PI uses an additional optimization to improve the policy via policy gradients. Therefore, PI requires an additional optimization compared to VI approaches.

---

### Algorithm 1 Continuous Fitted Value Iteration (cFVI)

---

**Input:** Dynamics Model  $f_c(\mathbf{x}, \mathbf{u})$  & Dataset  $\mathcal{D}$

**Result:** Value Function  $V^*(\mathbf{x}; \psi^*)$

---

**for**  $k = 0 \dots N$  **do**

    // Compute Value Target for  $\mathbf{x} \in \mathcal{D}$ :

$$V_{\text{tar}}(\mathbf{x}_i) = \int_0^T \beta \exp(-\beta t) R_t dt + \exp(-\beta T) R_T$$

$$R_t = \int_0^t \exp(-\rho\tau) r_c(\mathbf{x}_\tau, \mathbf{u}_\tau) d\tau + \exp(-\rho t) V^k(\mathbf{x}_t)$$

$$\mathbf{x}_\tau = \mathbf{x}_i + \int_0^\tau f_c(\mathbf{x}_t, \mathbf{u}_t) dt$$

$$\mathbf{u}_\tau = \nabla \tilde{g}(\mathbf{B}(\mathbf{x}_\tau) \nabla_x V^k(\mathbf{x}_\tau))$$

    // Fit Value Function:

$$\psi_{k+1} = \arg \min_{\psi} \sum_{\mathbf{x} \in \mathcal{D}} \|V_{\text{tar}}(\mathbf{x}) - V(\mathbf{x}; \psi)\|^p$$

**if** RTDP cFVI **then**

        // Add samples from  $\pi^{k+1}$  to FIFO buffer  $\mathcal{D}$

$$\mathcal{D}^{k+1} = h(\mathcal{D}^k, \{\mathbf{x}_0^{k+1} \dots \mathbf{x}_N^{k+1}\})$$

**end if**

**end for**

---

### 3.2. Deriving the Analytic Optimal Policy

One cannot directly apply FVI to continuous actions due to the maximization in Equation 8. To compute the value function target, one would need to solve an optimization in each iteration. To extend value iteration to continuous actions, we show that one can solve this maximization analytically for the considered continuous-time problem. This solution enables the efficient computation of the value function target and extends FVI to continuous actions.

**Theorem 1.** *If the dynamics are control affine (Equation 5), the reward is separable w.r.t. to state and action (Equation 6) and the action cost  $g_c$  is positive definite and strictly convex, the continuous-time optimal policy  $\pi^k$  is described by*

$$\pi^k(\mathbf{x}) = \nabla \tilde{g}_c(\mathbf{B}(\mathbf{x})^T \nabla_x V^k) \quad (10)$$

where  $\tilde{g}$  is the convex conjugate of  $g$  and  $\nabla_x V^k$  is the Jacobian of current value function  $V^k$  w.r.t. the system state.

*Proof Sketch.* The detailed proof is provided in the appendix. This derivation follows the work of Lutter et. al. (2019), which generalized the special cases initially described by Lyshevski (1998) and Doya (2000) to a wider class of reward functions. The value iteration target (Equation 8) can be expressed by using the Taylor series for the value function  $V(\mathbf{x}_{t+1})$  and substituting Equation 5 and Equation 6. This reformulated value target is described by

$$V_{\text{tar}} = \max_{\mathbf{u}} \left[ \gamma \nabla_x V^T (\mathbf{a} + \mathbf{B}\mathbf{u}) + \gamma O(\cdot) + g_c - g_c \right] \Delta t + \gamma V$$

with the higher order terms  $O(\Delta t, \mathbf{x}, \mathbf{u})$ . In the continuous-time limit the higher order terms disappear and  $\gamma = 1$ .

Therefore, the optimal action is described by

$$\mathbf{u}_t^* = \arg \max_{\mathbf{u}} \nabla_x V^T \mathbf{B}(\mathbf{x}_t) \mathbf{u} - g_c(\mathbf{u}). \quad (11)$$

Equation 11 can be solved analytically as  $g_c$  is strictly convex and hence  $\nabla g_c(\mathbf{u}) = \mathbf{w}$  is invertible, i.e.,  $\mathbf{u} = [\nabla g_c]^{-1}(\mathbf{w}) := \nabla \tilde{g}_c(\mathbf{w})$  with the convex conjugate  $\tilde{g}$ . The solution of Equation 11 is described by

$$\mathbf{B}^T \nabla_x V^* - \nabla g_c(\mathbf{u}) := 0 \Rightarrow \mathbf{u}^* = \nabla \tilde{g}_c \left( \mathbf{B}^T \nabla_x V^k \right).$$

□

Unrolling the closed form policy over time performs hill-climbing on the value function, where the step-size corresponds to the time-discretization. The inner part  $\mathbf{B}(\mathbf{x})^T \nabla_x V^k$  determines the direction of steepest ascent and the action cost rescales this direction. For example, a zero action cost for all admissible actions makes the optimal policy take the largest possible actions, i.e., the policy is bang-bang controller. A quadratic action cost linearly rescales the gradient. A barrier shaped cost clips the gradients. A complete guide on designing the action cost to shape the optimal policy can be found in Lutter et. al. (2019). The step-size, which corresponds to the control frequency of the discrete time controller, determines the convergence to the value function maximum. If the step size is too large the system becomes unstable and does not converge to the maximum. For most real world robotic systems with natural frequencies below 5Hz and control frequencies above 100Hz, the resulting step-size is sufficiently small to achieve convergence to the value function maximum. Therefore,  $\pi^*$  can be used for discrete time controllers. Furthermore, the continuous-time policy can be used for intermittent control (event-based control) where interacting with the system occurs at irregular time-steps and each interaction is associated with a cost (Aström, 2008).

### 3.3. Learning the Optimal Value Function

Using the closed-form policy that analytically solves the maximization, one can derive a computationally efficient value iteration update for continuous states and actions. Substituting  $\mathbf{u}^*$  (Equation 10) into the value iteration target computation (Equation 8) yields

$$V_{\text{tar}}(\mathbf{x}_t) = r \left( \mathbf{x}_t, \nabla \tilde{g} \left( \mathbf{B}(\mathbf{x}_t)^T \nabla_x V^k \right) \right) + \gamma V^k(\mathbf{x}_{t+1}; \psi_k)$$

with  $\mathbf{x}_{t+1} = f \left( \mathbf{x}_t, \nabla \tilde{g} \left( \mathbf{B}(\mathbf{x}_t)^T \nabla_x V^k \right) \right).$

Combined with fitting the value function to the target value (Equation 9), these two steps constitute the value function update of cFVI. Repeatedly computing the value target and fitting the value function leads to the optimal value function and policy.

For the continuous-time limit, the computation of the *naive* value function target should be adapted as the convergence

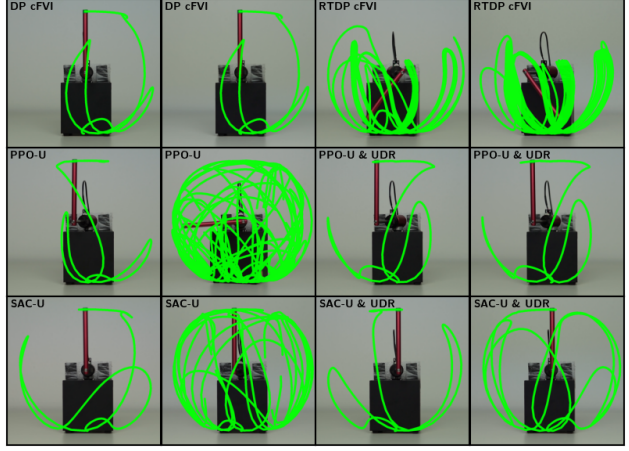


Figure 3. Tracked swing-up of the Furuta pendulum for best and worst roll out. DP cFVI can consistently swing-up the pendulum with minimal variance between roll outs. DP cFVI achieves much better performance compared to most baselines. Only PPO with domain randomization achieves comparable results. The remaining baselines are shown in the appendix.

decreases with decreasing time steps. As  $\Delta t$  decreases, the  $\gamma$  increases, i.e.,  $\gamma = \lim_{\Delta t \rightarrow 0} \exp(-\rho \Delta t) = 1$ . Therefore, the contraction coefficient of VI decreases exponentially with increasing sampling frequencies. This slower convergence is intuitive as higher sample frequencies effectively increase the number of steps to reach the goal.

**N-Step Value Function Target** To increase the computational efficiency cFVI, the exponentially weighted n-step value function target

$$V_{\text{tar}}(\mathbf{x}) = \int_0^T \beta \exp(-\beta t) R_t dt + \exp(-\beta T) R_T,$$

$$R_t = \int_0^t \exp(-\rho \tau) r_c(\mathbf{x}_\tau, \mathbf{u}_\tau) d\tau + \exp(-\rho t) V^k(\mathbf{x}_t),$$

and the exponential decay constant  $\beta$ , can be used. The integrals can be solved using any ordinary differential equation solver with fixed or adaptive step-size. We use the explicit Euler integrator with fixed steps to solve the integral for all samples in parallel using batched operations on the GPU. The nested integrals can be computed efficiently by recursively splitting the integral and reusing the estimate of the previous step. In practice we treat  $\beta$  as a hyperparameter and select  $T$  such that the weight of the  $R_T$  is  $\exp(-\beta T) := 10^{-4}$ .

The convergence rate improves as this target computation corresponds to the multi-step value target. The n-step target increases the contraction rate as this rate is proportional to  $\gamma^\ell$ . This approach is the continuous-time counterpart of the discrete eligibility trace of TD( $\lambda$ ) with  $\lambda = \exp(-\beta \Delta t)$  (Sutton et al., 1998). With respect to deep RL this discounted n-step value target is similar to the generalized advantage estimation (GAE) of PPO (Schulman et al., 2015; 2017) and model-based value expansion (MVE) (Feinberg et al., 2018;

Buckman et al., 2018). GAE and MVE have shown that the  $n$ -step target increases the sample efficiency and lead to faster convergence to the optimal policy.

**Offline and Online cFVI** The proposed approach is off-policy as the samples in the replay memory do not need to originate from the current policy  $\pi_k$ . Therefore, the dataset can either consist of a fixed dataset (i.e., batch or offline RL) or be updated within each iteration. In the offline dynamic programming case, the dataset contains samples from the compact state domain  $\mathcal{X}$ . We refer to the offline variant as DP cFVI. In the online case, the replay memory is updated with samples generated by the current policy  $\pi_k$ . Every iteration the states of the previous  $n$ -rollouts are added to the data and replace the oldest samples. This online update of state distribution performs real-time dynamic programming (RTDP) (Barto et al., 1995). We refer to the online variant as RTDP cFVI. The pseudo code of DP cFVI and RTDP cFVI is summarized in Algorithm 1.

#### 4. Value Function Hypothesis Space

The previous sections focused on learning the optimal value function independent of the value function representation. In this section, we focus on the value function representation.

**Value Function Representation** Most recent approaches use a fully connected network to approximate the value function. However, for the many tasks the hypothesis space of the value function can be narrowed. For control tasks, the state cost is often a negative distance measure between  $\mathbf{x}_t$  and the desired state  $\mathbf{x}_{\text{des}}$ . Hence,  $q_c$  is negative definite, i.e.,  $q(\mathbf{x}) < 0 \forall \mathbf{x} \neq \mathbf{x}_{\text{des}}$  and  $q(\mathbf{x}_{\text{des}}) = 0$ . These properties imply that  $V^*$  is a negative Lyapunov function, as  $V^*$  is negative definite,  $V^*(\mathbf{x}_{\text{des}}) = 0$  and  $\nabla_{\mathbf{x}} V^*(\mathbf{x}_{\text{des}}) = \mathbf{0}$  (Khalil & Grizzle, 2002). With a deep network a similar representation can be achieved by

$$V(\mathbf{x}; \psi) = -(\mathbf{x} - \mathbf{x}_{\text{des}})^T \mathbf{L}(\mathbf{x}; \psi) \mathbf{L}(\mathbf{x}; \psi)^T (\mathbf{x} - \mathbf{x}_{\text{des}})$$

with  $\mathbf{L}$  being a lower triangular matrix with positive diagonal. This positive diagonal ensures that  $\mathbf{L}\mathbf{L}^T$  is positive definite. Simply applying a ReLU activation to the last layer of a deep network is not sufficient as this would also zero the actions for the positive values and  $\nabla_{\mathbf{x}} V^*(\mathbf{x}_{\text{des}}) = \mathbf{0}$  cannot be guaranteed. The local quadratic representation guarantees that the gradient and hence, the action, is zero at the desired state. However, this representation can also not guarantee that the value function has only a single extrema at  $\mathbf{x}_{\text{des}}$  as required by the Lyapunov theory. In practice, the local regularization of the quadratic structure to avoid high curvature approximations is sufficient as the global structure is defined by the value function target.  $\mathbf{L}$  is the mean of a deep network ensemble with  $N$  independent parameters  $\psi_i$ . The ensemble mean smoothes the initial value function and is differentiable. Similar representations have been used by prior works in the safe reinforcement learning community

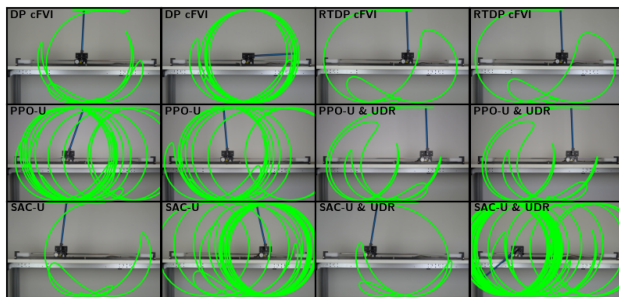


Figure 4. Tracked swing-up of the cartpole for the best and worst roll out. DP cFVI and RTDP cFVI can achieve the task. In the failure case, cFVI excels as the policy remains in the center. In contrast the deep RL baselines move the cart between the joint limits in case of failure. All baselines are shown in the appendix.

(Berkenkamp et al., 2017; Richards et al., 2018; Kolter & Manek, 2019; Chang et al., 2019; Bharadhwaj et al., 2021).

**Gradient Projection of State Transformations** State transformations should be explicitly incorporated in the value function and should not be contained in the environment, as for example in the openAI Gym (Brockman et al., 2016). If the transformation is implicit,  $\nabla_{\mathbf{x}} V$  might not be sensible. For example, the standard feature transform for a continuous revolute joint maps the joint state  $\mathbf{x} = [\theta, \dot{\theta}]$  to  $\mathbf{z} = [\sin(\theta), \cos(\theta), \dot{\theta}]$ . In this case the transformation  $h(\mathbf{x})$  must be included in the value function as the transformed state  $\mathbf{z}$  is on the tube shaped manifold. Hence, the naive gradient of  $V$  might not be in the manifold tangent space.

If the state transform is explicitly incorporated in the value function, this problem does not occur. The transformation can be included explicitly by  $V(\mathbf{x}; \psi) = f(h(\mathbf{x}); \psi)$  and  $\nabla_{\mathbf{x}} V(\mathbf{x}; \psi) = \partial f(h(\mathbf{x}); \psi) / \partial h \partial h(\mathbf{x}) / \partial \mathbf{x}$ . In this case, the gradient of the transformed state is projected to the tangent space of the feature transform. Therefore, the value function gradient points in a plausible direction.

## 5. Experiments

In the following the experimental setup and results are described. The exact experiment specification, all qualitative plots and an additional ablation study on model ensembles is provided in the appendix.

### 5.1. Experimental Setup

**Systems** The algorithms are compared using the standard non-linear control benchmark, the *swing-up* of under-actuated systems. Specifically, we apply the algorithms to the torque-limited pendulum, cartpole (Figure 4) and Furuta pendulum (Figure 3). For the physical systems the dynamics model of the manufacturer is used (Quanser, 2018). The control frequency is optimized for each algorithm. The task is considered solved, if the pendulum angle  $\alpha$  is below  $\pm 5^\circ$  degree for the last second.

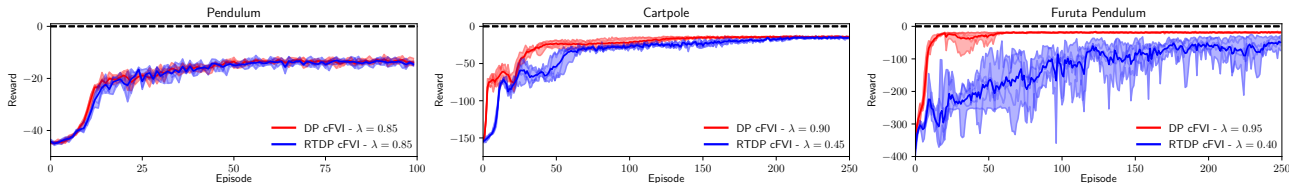


Figure 5. The learning curves for DP cFVI and RTDP cFVI averaged over 5 seeds. The shaded area displays the *min/max* range between seeds. DP cFVI achieves consistent learning of the optimal policy with very low variance between seeds. RTDP cFVI has a higher variation and learns slower compared to DP cFVI. RTDP cFVI needs to discover the steady-state state distribution first. Especially for the Furuta pendulum the range increases between seeds as minor changes in the policy lead to large changes of state-distribution.

Table 1. Average rewards on the simulated and physical systems. The ranking describes the decrease in reward compared to the best result averaged on all systems. The initial state distribution during training is noted by  $\mu$ . The dynamics are either deterministic model  $\theta \sim \delta(\theta)$  or sampled using uniform domain randomization  $\theta \sim \mathcal{U}(\theta)$ . During evaluation the roll outs start with the pendulum pointing downwards. DP cFVI obtains high rewards on all systems and is the highest ranking algorithm compared to the baselines.

Algorithm	$\mu$	$\theta$	Simulated Pendulum		Simulated Cartpole		Simulated Furuta Pendulum		Physical Cartpole		Physical Furuta Pendulum		Average Ranking [%]
			Success [%]	Reward [ $\mu \pm 1.96\sigma$ ]	Success [%]	Reward [ $\mu \pm 1.96\sigma$ ]	Success [%]	Reward [ $\mu \pm 1.96\sigma$ ]	Success [%]	Reward [ $\mu \pm 1.96\sigma$ ]	Success [%]	Reward [ $\mu \pm 1.96\sigma$ ]	
DP cFVI (ours)	$\mathcal{U}$	$\mathcal{U}(\theta)$	100.0	<b>-030.5 <math>\pm</math> 000.8</b>	100.0	<b>-024.2 <math>\pm</math> 002.1</b>	100.0	<b>-027.7 <math>\pm</math> 001.6</b>	73.3	-143.7 $\pm$ 210.4	100.0	<b>-082.1 <math>\pm</math> 007.6</b>	<b>-008.8</b>
RTDP cFVI (ours)	$\mathcal{U}$	$\mathcal{U}(\theta)$	100.0	<b>-031.1 <math>\pm</math> 001.4</b>	100.0	<b>-024.9 <math>\pm</math> 001.6</b>	100.0	-040.1 $\pm$ 002.7	100.0	<b>-101.1 <math>\pm</math> 029.0</b>	00.0	-1009.9 $\pm$ 004.5	-240.4
SAC	$\mathcal{N}$	$\mathcal{U}(\theta)$	100.0	<b>-031.1 <math>\pm</math> 000.1</b>	100.0	-026.9 $\pm$ 003.2	100.0	-029.3 $\pm$ 001.5	00.0	-518.6 $\pm$ 028.1	86.7	-330.7 $\pm$ 799.0	-148.4
SAC & UDR	$\mathcal{N}$	$\delta(\theta)$	100.0	-032.9 $\pm$ 000.6	100.0	-029.7 $\pm$ 004.6	100.0	-032.0 $\pm$ 001.1	100.0	-394.8 $\pm$ 382.8	100.0	-181.4 $\pm$ 157.9	-092.3
SAC	$\mathcal{U}$	$\mathcal{U}(\theta)$	100.0	<b>-030.6 <math>\pm</math> 001.4</b>	100.0	<b>-024.2 <math>\pm</math> 001.4</b>	100.0	<b>-028.1 <math>\pm</math> 002.0</b>	53.3	-144.5 $\pm$ 204.0	100.0	-350.8 $\pm$ 433.3	-076.1
SAC & UDR	$\mathcal{U}$	$\mathcal{U}(\theta)$	100.0	-031.4 $\pm$ 002.5	100.0	<b>-024.2 <math>\pm</math> 001.3</b>	100.0	<b>-028.1 <math>\pm</math> 001.3</b>	40.0	-296.4 $\pm$ 418.9	100.0	-092.3 $\pm$ 064.1	-042.4
DDPG	$\mathcal{N}$	$\mathcal{U}(\theta)$	100.0	<b>-031.1 <math>\pm</math> 000.4</b>	98.0	-050.4 $\pm$ 285.6	100.0	-030.5 $\pm$ 003.5	06.7	-536.7 $\pm$ 262.7	46.7	-614.1 $\pm$ 597.8	-242.6
DDPG & UDR	$\mathcal{N}$	$\delta(\theta)$	100.0	-032.5 $\pm$ 000.5	100.0	-027.4 $\pm$ 002.3	100.0	-034.6 $\pm$ 009.8	00.0	-517.9 $\pm$ 117.6	86.7	-192.7 $\pm$ 404.8	-119.3
DDPG	$\mathcal{U}$	$\mathcal{U}(\theta)$	100.0	-031.5 $\pm$ 000.7	100.0	-028.2 $\pm$ 005.5	100.0	-030.0 $\pm$ 001.7	06.7	-459.4 $\pm$ 248.3	100.0	-146.6 $\pm$ 218.3	-092.9
DDPG & UDR	$\mathcal{U}$	$\mathcal{U}(\theta)$	100.0	-032.5 $\pm$ 003.6	100.0	-027.2 $\pm$ 001.0	100.0	-032.1 $\pm$ 001.5	00.0	-318.1 $\pm$ 063.4	100.0	-156.7 $\pm$ 246.4	-068.8
PPO	$\mathcal{N}$	$\mathcal{U}(\theta)$	100.0	-032.0 $\pm$ 000.2	100.0	-031.5 $\pm$ 007.2	100.0	-081.1 $\pm$ 018.3	00.0	-287.9 $\pm$ 068.8	33.3	-718.7 $\pm$ 456.1	-240.9
PPO & UDR	$\mathcal{N}$	$\delta(\theta)$	100.0	-032.3 $\pm$ 000.6	100.0	-084.0 $\pm$ 007.8	100.0	-040.9 $\pm$ 004.6	00.0	-435.4 $\pm$ 111.9	46.7	-935.7 $\pm$ 711.6	-338.6
PPO	$\mathcal{U}$	$\mathcal{U}(\theta)$	100.0	-033.4 $\pm$ 004.7	99.0	-039.7 $\pm$ 045.7	100.0	-038.2 $\pm$ 013.1	00.0	-183.8 $\pm$ 018.0	60.0	-755.3 $\pm$ 811.0	-206.1
PPO & UDR	$\mathcal{U}$	$\mathcal{U}(\theta)$	100.0	-035.6 $\pm$ 003.1	100.0	-044.8 $\pm$ 021.4	100.0	-048.5 $\pm$ 006.2	40.0	-143.8 $\pm$ 016.1	100.0	<b>-080.6 <math>\pm</math> 010.8</b>	-044.0

**Baselines** The performance is compared to the actor-critic deep RL methods: DDPG (Lillicrap et al., 2015), SAC (Haarnoja et al., 2018) and PPO (Schulman et al., 2017). We compare two different initial state distributions  $\mu$ . For  $\mu = \mathcal{U}$ , the initial pendulum angle  $\alpha_0$  is sampled uniformly  $\alpha_0 \sim \mathcal{U}(-\pi, +\pi)$ . For  $\mu = \mathcal{N}$ , the initial angle is sampled from a Gaussian distribution with the pendulum facing downwards  $\alpha_0 \sim \mathcal{N}(\pm\pi, \sigma)$ . The uniform sampling avoids the exploration problem and generates a larger state distribution of the optimal policy. In addition we augment each baseline with uniform domain randomization (Muratore et al., 2018).

## 5.2. Simulation Results

The average rewards of our method - Continuous Fitted Value Iteration (cFVI) and the baselines are summarized in Table 1. The learning curves are shown in Figure 5. In simulation, DP cFVI, the offline version of cFVI with fixed dataset, achieves the best rewards for all three systems and marginally outperforms SAC in terms of average reward. It’s important to point out that identical mean rewards do not imply a similar qualitative performance. For example, DP cFVI and SAC-U achieve identical mean reward on the cartpole despite having very different closed-loop dynamics (Figure 11 - Appendix). Notably, RTDP cFVI, which uses the state distribution of the current policy rather than a fixed dataset, solves the tasks for all three systems. For the

pendulum and the cartpole, the reward is comparable to the best performing algorithms. While for the Furuta pendulum, the reward is lower (higher) than the best (worst) performing Deep RL baselines.

The learning curves highlight that the variance between seeds is very low for DP cFVI. It is important to note that Figure 5 shows the min-max range of the different seeds rather than the confidence intervals or the standard error. Therefore, the exact weight initialization and sampling of the fixed dataset, which is different for each seed, do not have a large impact on the learning progress or obtained solution. For RTDP cFVI the variance between seed increases and the learning speed decreases compared to DP cFVI. The pendulum is an outlier for RTDP cFVI, as the state-distribution is nearly independent of the current policy if the pendulum angle is initialized uniformly. Especially for the Furuta pendulum the variance increases and learning speed decreases. For this system slightly different policies can cause vastly different state distributions due to the low condition number of the mass-matrix. Therefore, RTDP cFVI takes longer to reach a stationary distribution and increase the obtained reward.

In general it is important to point out that the exploration for RTDP cFVI is more challenging compared to other deep RL methods as RTDP cFVI uses a higher control frequency. Due to the shorter steps simple Gaussian noise

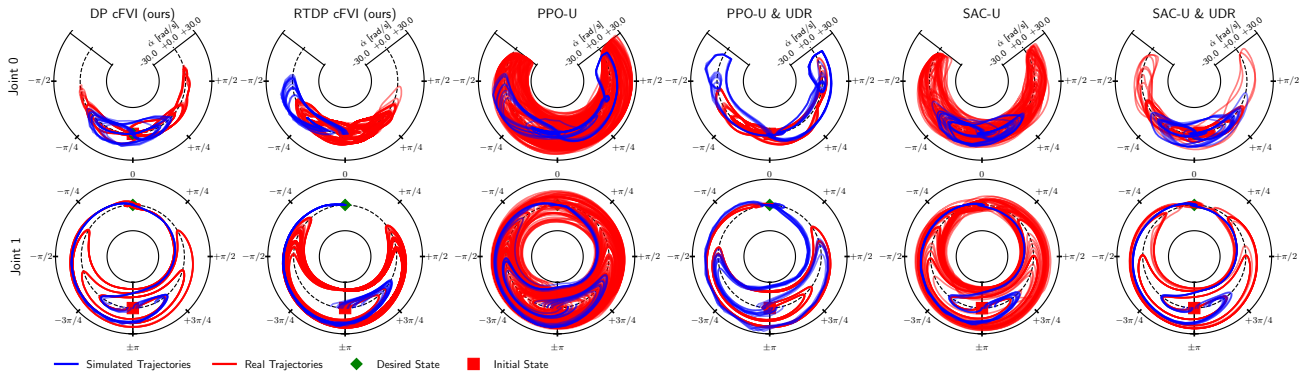


Figure 6. The simulated (blue) and real world (red) roll outs for the Furuta pendulum controlled by cFVI and the baselines. The deviation from the dashed line denotes the joint velocity. The distribution shift is clearly visible as the blue and red trajectories do not overlap, e.g., cFVI requires more pre-swings on the physical system. For PPO-U the distribution shift causes random exploration of the complete state-domain. cFVI achieves the best qualitative performance as it can swing-up the pendulum from both directions and achieves nearly identical roll-outs. In contrast the deep RL baselines have higher distribution shift due to the dynamics mismatch. The trajectories are also less consistent and vary significantly between roll outs. The plots for all baselines and the physical cartpole are shown in the appendix.

averages out and does not lead to a large exploration in state space. Therefore, the performance and variance of RTDP cFVI could be improved using a more coherent exploration strategy in future work.

### 5.3. Ablation Study - $N$ -step Value Targets

The learning curves, averaged over 5 seeds, for different  $n$ -step value targets are shown in Figure 7. When increasing  $\lambda$ , which implicitly increases the  $n$ -step horizon (section 3.3), the convergence rate to the optimal value function increases. This increased learning speed is expected as Equation 7 shows that the convergence rate depends on  $\gamma^n$  with  $\gamma < 1$ . While the learning speed measured in iterations increases with  $\lambda$ , the computational complexity also increases. Longer horizons require to simulate  $n$  sequential steps increasing the required wall-clock time. Therefore, the computation time increases exponentially with increasing  $\lambda$ . For example the forward roll out in every iteration of the pendulum increases exponentially from 0.4s ( $\lambda = 0.01$ ) to 56.4s ( $\lambda = 0.99$ )<sup>1</sup>. For the Furuta pendulum and the cartpole extremely long horizons of 100+ steps start to diverge as the value function target over fits to the untrained value function. For RTDP cFVI, the horizons must smaller, i.e., 10 - 20 steps. For longer horizons the predicted rollout overfits to the value function outside of the current state distribution, which prevents learning or leads to pre-mature convergence (see Appendix Figure 9). Therefore, DP cFVI works best with  $\lambda \in [0.85, 0.95]$  and RTDP cFVI with  $\lambda \in [0.45 - 0.55]$ .

### 5.4. Physical Experiments

The quantitative results of the physical experiments are summarized in Table 1<sup>2</sup>, with additional plots and baselines in the appendix. DP cFVI achieves high reward on both

physical systems and outperforms the baselines. Only PPO-U UDR achieves comparable performance on both physical systems. On average the performance of the deep RL variants increases with with larger initial state distribution.

The trajectories of the Furuta Pendulum are shown in Figure 6. The distribution shift between the physical and simulated trajectories is clearly visible as the trajectories largely deviate. DP cFVI achieves a highly repeatable performance on the noisy physical system. All 15 roll-outs have a similar trajectory. The baselines require domain randomization and uniform initial distribution to achieve a reward comparable to DP cFVI. All other deep RL baselines have a high variance between roll-outs. For example the simulation gap causes PPO-U to randomly explore the complete state space. RTDP cFVI does not achieve the swing-up but fails gracefully as it stabilizes the pendulum on a stable limit cycle and does not hit the joint limits. Furthermore, this limit cycle is highly repetitive and the variance in reward is low. For many robotic tasks, it is preferable to fail gracefully rather than hitting the limits and solving the task.

On the physical cartpole, DP cFVI obtains a higher reward than RTDP cFVI, if the swing-up is successful. However, RTDP cFVI has a higher average reward across all roll-outs as DP cFVI only achieves a 73% success rate. If DP cFVI fails, the resulting behavior is a deterministic limit-cycle where the cart is centered. The policy decelerates the pendulum, but due to the backlash in the linear actuator the applied force is not sufficient. Most deep RL baselines do not achieve a consistent swing-up. Only PPO-U UDR, SAC-U achieve a repeatable high reward. SAC-N UDR achieves the swing-up but always hits the joint limit and obtains a relative low reward. The failure cases of all deep RL baselines are stochastic and repeatedly hit the joint limit. If the pendulum over-swings the cart starts to oscillate

<sup>1</sup>Wall-clock time on an AMD 3900X and a NVIDIA RTX 3090

<sup>2</sup>Task Video available at: <https://sites.google.com/view/value-iteration>

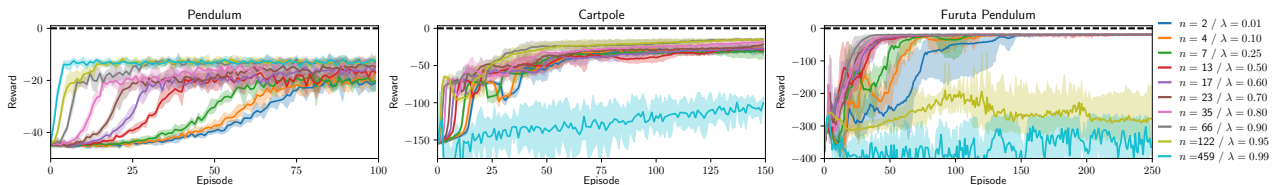


Figure 7. The learning curves averaged over 5 seeds for the  $n$ -step value function target. The shaded area displays the  $\min/\max$  range between seeds. The step count is selected such that  $\lambda^n := 10^{-4}$ . Increasing the horizon of the value function target increases the convergence rate to the optimal value function. For very long horizons the learning diverges as it over fits to the current value function approximation. Furthermore, the performance of the optimal policy also increases with roll out length.

between both joint limits. Notably, domain randomization does not improve the performance on the physical cartpole. For this system, the simulation gap originates from backlash and friction in the actuator rather than the uncertainty of the model parameters. As the actuation is not simulated, these parameters cannot be randomized. For the cartpole, the larger initial state distribution is required to achieve high rewards on the physical system. All baselines with a gaussian initial state distribution achieve only a low reward. See appendix for qualitative simulated and physical results.

## 6. Related Work

Continuous-time RL started with the seminal work of Doya (2000). Since then, various approaches have been proposed to solve the Hamilton-Jacobi-Bellman (HJB) differential equation with the machine learning toolset. These methods can be divided into trajectory and state-space based methods.

Trajectory based methods solve the stochastic HJB along a trajectory to obtain the optimal trajectory. For example path integral control uses the non-linear, control-affine dynamics with quadratic action costs to simplify the HJB to a linear partial differential equation (Kappen, 2005; Todorov, 2007; Theodorou et al., 2010; Pan et al., 2014). This differential equation can be transformed to a path-integral using the Feynman-Kac formulae. The path-integral can then be solved using Monte Carlo sampling to obtain the optimal state and action sequence. Recently, this approach has been used in combination with deep networks (Rajagopal et al., 2016; Pereira et al., 2019; 2020).

State-space based methods solve the HJB globally to obtain an optimal non-linear controller applicable on the complete state domain. Classical approaches discretize the continuous spaces into a grid and solve the HJB or the robust Hamilton-Jacobi-Isaac (HJI) using a PDE solver (Bansal et al., 2017). To overcome the curse of dimensionality of the grid based methods, machine learning methods that use function approximation and sampling have been proposed. For example, regression based approaches solved the HJB by fitting a radial-basis-function networks (Doya, 2000), deep networks (Tassa & Erez, 2007; Lutter et al., 2019; Kim et al., 2020), kernels (Hennig, 2011) or polynomial functions (Yang et al., 2014; Liu et al., 2014) to minimize

the HJB residual. The naive optimization of this objective, while successful for different PDEs (Raissi et al., 2017a;b; 2020), does not work for the HJB as the exact location of the boundary condition is unknown. Therefore, authors used various optimization tricks such as annealing the noise (Tassa & Erez, 2007) or the discounting (Lutter et al., 2019) to obtain the optimal value function.

In this work we propose a DP based algorithm. Instead of trying to solve the HJB PDE directly via regression as previous learning-based methods, we leverage the continuous-time domain and solve the HJB by iteratively applying the Bellman optimality principle. The resulting algorithm has a simpler and more robust optimization compared to the previous approaches using direct regression.

## 7. Conclusion and Future Work

We proposed continuous fitted value iteration (cFVI). This algorithm enables dynamic programming with known model for problems with continuous states and actions without discretization. Therefore, cFVI avoids the curse of dimensionality associated with discretization and the policy optimization of actor-critic approaches. Exploiting the insights from the continuous-time formulation, the optimal action can be computed analytically. This closed-form solution permits the efficient computation of the value function target and enables the extension of value iteration. The non-linear control experiments showed that value iteration on the compact state space has the same quantitative performance as deep RL methods in simulation. On the sim2real tasks, DP cFVI excels compared deep RL algorithms that include domain randomization and use uniform initial state distribution.

In future work, we plan to extend cFVI to (1) offline/batch model-based RL, (2) stochastic maximum entropy policies and (3) explicit model robustness. cFVI uses a fixed data set and hence, can be combined with learning control-affine models to the offline RL benchmark (Gulcehre et al., 2020; Mandlekar et al., 2019). We plan to extend cFVI to maximum entropy and stochastic policies to improve the exploration of RTDP cFVI. Finally, we plan to incorporate robustness w.r.t. to changes in dynamics by using the adversarial continuous RL formulation culminating in the Hamilton-Jacobi-Isaacs differential equation rather than the HJB (Isaacs, 1999).



## Acknowledgements

M. Lutter was an intern at Nvidia during this project. A. Garg was partially supported by CIFAR AI Chair. We also want to thank Fabio Muratore, Joe Watson and the ICML reviewers for their feedback. Furthermore, we want to thank the open-source projects SimuRLacra (Muratore, 2020), MushroomRL (D’Eramo et al., 2020), NumPy (Harris et al., 2020) and PyTorch (Paszke et al., 2019).

## References

- Aström, K. J. Event based control. In *Analysis and design of nonlinear control systems*, pp. 127–147. Springer, 2008.
- Baird, L. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings 1995*, pp. 30–37. Elsevier, 1995.
- Bansal, S., Chen, M., Herbert, S., and Tomlin, C. J. Hamilton-Jacobi reachability: A brief overview and recent advances. 2017.
- Barto, A. G., Bradtke, S. J., and Singh, S. P. Learning to act using real-time dynamic programming. *Artificial intelligence*, 72(1-2):81–138, 1995.
- Bellman, R. *Dynamic Programming*. Princeton University Press, USA, 1957. ISBN 0691146683.
- Berkenkamp, F., Turchetta, M., Schoellig, A., and Krause, A. Safe model-based reinforcement learning with stability guarantees. In *Advances in neural information processing systems*, pp. 908–918, 2017.
- Bharadhwaj, H., Kumar, A., Rhinehart, N., Levine, S., Shkurti, F., and Garg, A. Conservative safety critics for exploration. In *International Conference on Learning Representations (ICLR)*, 2021.
- Boyan, J. and Moore, A. Generalization in reinforcement learning: Safely approximating the value function. *Advances in neural information processing systems*, 7:369–376, 1994.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- Buckman, J., Hafner, D., Tucker, G., Brevdo, E., and Lee, H. Sample-efficient reinforcement learning with stochastic ensemble value expansion. In *Advances in Neural Information Processing Systems*, pp. 8224–8234, 2018.
- Chang, Y.-C., Roohi, N., and Gao, S. Neural lyapunov control. In *Advances in Neural Information Processing Systems*, pp. 3245–3254, 2019.
- Chebatar, Y., Handa, A., Makoviyuchuk, V., Macklin, M., Issac, J., Ratliff, N., and Fox, D. Closing the sim-to-real loop: Adapting simulation randomization with real world experience, 2019.
- Da, X., Xie, Z., Hoeller, D., Boots, B., Anandkumar, A., Zhu, Y., Babich, B., and Garg, A. Learning a Contact-Adaptive Controller for Robust, Efficient Legged Locomotion. In *Conference on Robot Learning (CoRL)*, 2020.
- D’Eramo, C., Tateo, D., Bonarini, A., Restelli, M., and Peters, J. Mushroomrl: Simplifying reinforcement learning research. <https://github.com/MushroomRL/mushroom-rl>, 2020.
- Doya, K. Reinforcement learning in continuous time and space. *Neural computation*, 12(1):219–245, 2000.
- Ernst, D., Geurts, P., and Wehenkel, L. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6(Apr):503–556, 2005.
- Farahmand, A. M., Ghavamzadeh, M., Szepesvári, C., and Mannor, S. Regularized fitted q-iteration for planning in continuous-space markovian decision problems. In *2009 American Control Conference*, pp. 725–730. IEEE, 2009.
- Feinberg, V., Wan, A., Stoica, I., Jordan, M. I., Gonzalez, J. E., and Levine, S. Model-based value estimation for efficient model-free reinforcement learning. *arXiv preprint arXiv:1803.00101*, 2018.
- Gulcehre, C., Wang, Z., Novikov, A., Paine, T. L., Colmenarejo, S. G., Zolna, K., Agarwal, R., Merel, J., Mankowitz, D., Paduraru, C., et al. RL unplugged: Benchmarks for offline reinforcement learning. *arXiv preprint arXiv:2006.13888*, 2020.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*, 2018.
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. Array programming with NumPy. *Nature*, 2020.
- Harrison\*, J., Garg\*, A., Ivanovic, B., Zhu, Y., Savarese, S., Fei-Fei, L., and Pavone (\* equal contribution), M. AdaPT: Zero-Shot Adaptive Policy Transfer for Stochastic Dynamical Systems. In *International Symposium on Robotics Research (ISRR)*. Springer STAR, 2017.

- Hennig, P. Optimal reinforcement learning for gaussian systems. In *Advances in Neural Information Processing Systems*, pp. 325–333, 2011.
- Isaacs, R. *Differential games: a mathematical theory with applications to warfare and pursuit, control and optimization*. Courier Corporation, 1999.
- Kappen, H. J. Linear theory for control of nonlinear stochastic systems. *Physical review letters*, 95(20):200201, 2005.
- Khalil, H. K. and Grizzle, J. W. *Nonlinear systems*, volume 3. Prentice hall Upper Saddle River, NJ, 2002.
- Kim, J., Shin, J., and Yang, I. Hamilton-jacobi deep q-learning for deterministic continuous-time systems with lipschitz continuous controls. *arXiv preprint arXiv:2010.14087*, 2020.
- Kirk, D. E. *Optimal control theory: an introduction*. Courier Corporation, 1970.
- Kolter, J. Z. and Manek, G. Learning stable deep dynamics models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Liu, D., Wang, D., Wang, F.-Y., Li, H., and Yang, X. Neural-network-based online hjb solution for optimal robust guaranteed cost control of continuous-time uncertain nonlinear systems. *IEEE transactions on cybernetics*, 44(12):2834–2847, 2014.
- Lutter, M., Belousov, B., Listmann, K., Clever, D., and Peters, J. HJB optimal feedback control with deep differential value functions and action constraints. In *Conference on Robot Learning (CoRL)*, 2019.
- Lyshevski, S. E. Optimal control of nonlinear continuous-time systems: design of bounded controllers via generalized nonquadratic functionals. In *Proceedings of the 1998 American Control Conference*, volume 1, pp. 205–209. IEEE, 1998.
- Mandlekar\*, A., Zhu\*, Y., Garg\*, A., Fei-Fei, L., and Savarese (\* equal contribution), S. Adversarially Robust Policy Learning through Active Construction of Physically-Plausible Perturbations. In *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2017.
- Mandlekar, A., Booher, J., Spero, M., Tung, A., Gupta, A., Zhu, Y., Garg, A., Savarese, S., and Fei-Fei, L. Scaling robot supervision to hundreds of hours with roboturk: Robotic manipulation dataset through human reasoning and dexterity. In *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2019.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Hiedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529, 2015.
- Munos, R. and Szepesvári, C. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research*, 9 (May):815–857, 2008.
- Muratore, F. Simurlacra - a framework for reinforcement learning from randomized simulations. <https://github.com/famura/SimuRLacra>, 2020.
- Muratore, F., Treede, F., Gienger, M., and Peters, J. Domain randomization for simulation-based policy optimization with transferability assessment. In *Conference on Robot Learning*, pp. 700–713, 2018.
- OpenAI, Andrychowicz, M., Baker, B., Chociej, M., Jozefowicz, R., McGrew, B., Pachocki, J., Petron, A., Plappert, M., Powell, G., Ray, A., Schneider, J., Sidor, S., Tobin, J., Welinder, P., Weng, L., and Zaremba, W. Learning dexterous in-hand manipulation, 2019.
- Pan, Y., Theodorou, E. A., and Kontitsis, M. Model-based path integral stochastic control: A bayesian nonparametric approach. *arXiv preprint arXiv:1412.3038*, 2014.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*. 2019.
- Pereira, M., Wang, Z., Exarchos, I., and Theodorou, E. Learning deep stochastic optimal control policies using forward-backward sdes. In *Robotics: science and systems*, 2019.
- Pereira, M. A., Wang, Z., Exarchos, I., and Theodorou, E. A. Safe optimal control using stochastic barrier functions and deep forward-backward sdes. *arXiv preprint arXiv:2009.01196*, 2020.
- Puterman, M. L. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 1994.
- Quanser. Quanser courseware and resources. <https://www.quanser.com/solution/control-systems/>, 2018.

- Raissi, M., Perdikaris, P., and Karniadakis, G. E. Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. *arXiv preprint arXiv:1711.10561*, 2017a.
- Raissi, M., Perdikaris, P., and Karniadakis, G. E. Physics informed deep learning (part ii): data-driven discovery of nonlinear partial differential equations. *arXiv preprint arXiv:1711.10566*, 2017b.
- Raissi, M., Yazdani, A., and Karniadakis, G. E. Hidden fluid mechanics: Learning velocity and pressure fields from flow visualizations. *Science*, 367(6481):1026–1030, 2020.
- Rajagopal, K., Balakrishnan, S. N., and Busemeyer, J. R. Neural network-based solutions for stochastic optimal control using path integrals. *IEEE transactions on neural networks and learning systems*, 28(3):534–545, 2016.
- Ramos, F., Possas, R. C., and Fox, D. Bayessim: adaptive domain randomization via probabilistic inference for robotics simulators, 2019.
- Richards, S. M., Berkenkamp, F., and Krause, A. The lyapunov neural network: Adaptive stability certification for safe learning of dynamical systems. *arXiv preprint arXiv:1808.00924*, 2018.
- Riedmiller, M. Neural fitted q iteration—first experiences with a data efficient neural reinforcement learning method. In *European Conference on Machine Learning*, pp. 317–328. Springer, 2005.
- Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.
- Sutton, R. S., Barto, A. G., et al. *Introduction to reinforcement learning*, volume 135. MIT press Cambridge, 1998.
- Tassa, Y. and Erez, T. Least squares solutions of the hjb equation with neural network value-function approximators. *IEEE transactions on neural networks*, 18(4):1031–1041, 2007.
- Tesauro, G. Practical issues in temporal difference learning. *Machine learning*, 8(3):257–277, 1992.
- Theodorou, E., Buchli, J., and Schaal, S. Reinforcement learning of motor skills in high dimensions: A path integral approach. In *2010 IEEE International Conference on Robotics and Automation*, pp. 2397–2403. IEEE, 2010.
- Todorov, E. Linearly-solvable markov decision problems. In *Advances in neural information processing systems*, pp. 1369–1376, 2007.
- Tsitsiklis, J. N. and Van Roy, B. Feature-based methods for large scale dynamic programming. *Machine Learning*, 22(1-3):59–94, 1996.
- Xie, Z., Da, X., van de Panne, M., Babich, B., and Garg, A. Dynamics Randomization Revisited: A Case Study for Quadrupedal Locomotion. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2021.
- Yang, X., Liu, D., and Wang, D. Reinforcement learning for adaptive optimal control of unknown continuous-time nonlinear systems with input constraints. *International Journal of Control*, 87(3):553–566, 2014.

## 8. Appendix

### Proof of Theorem 1

**Theorem.** *If the dynamics are control affine (Equation 5), the reward is separable w.r.t. to state and action (Equation 6) and the action cost  $g_c$  is positive definite and strictly convex, the continuous time optimal policy  $\pi^k$  w.r.t.  $V^k$  is described by*

$$\pi^k(\mathbf{x}) = \nabla \tilde{g}_c \left( \mathbf{B}(\mathbf{x})^T \nabla_x V^k \right) \quad (12)$$

with the convex conjugate  $\tilde{g}$  of  $g$  and the Jacobian of  $V$  w.r.t. the system state  $\nabla_x V$ .

*Proof.* This proof follows the derivation of Lutter et. al. (2019). This prior work derived the optimal policy  $\pi^*$  using the Hamilton Jacobi Bellman differential equation (HJB) and generalized the special case described by Doya (2000). The value iteration update (Equation 8) is defined as

$$V_{\text{tar}}(\mathbf{x}_t) = \max_{\mathbf{u}} r(\mathbf{x}_t, \mathbf{u}) + \gamma V(f(\mathbf{x}_t, \mathbf{u}); \psi_k).$$

Substituting the assumptions and using the Taylor expansion of the Value function, i.e.,

$$V(f(\mathbf{x}_t, \mathbf{u})) = V(\mathbf{x}_t) + \nabla_x V^T f_c(\mathbf{x}_t, \mathbf{u}) \Delta t + O(\Delta t, \mathbf{x}_t, \mathbf{u}) \Delta t,$$

this update can be rewritten - omitting all functional dependencies for brevity - as

$$\begin{aligned} V_{\text{tar}} &= \max_{\mathbf{u}} r + \gamma V + \gamma \nabla_x V^T f_c \Delta t + \gamma O \Delta t \\ &= \max_{\mathbf{u}} \left[ \gamma \nabla_x V^T (\mathbf{a} + \mathbf{B}\mathbf{u}) + \gamma O - g_c \right] \Delta t + \gamma V + q_c \Delta t \end{aligned}$$

with the higher order Terms  $O(\Delta t, \mathbf{x}_t, \mathbf{u})$ . Therefore, the optimal action is defined as

$$\mathbf{u}_t^* = \arg \max_{\mathbf{u}} \gamma \nabla_x V^T (\mathbf{a} + \mathbf{B}\mathbf{u}) + \gamma O(\Delta t, \mathbf{x}_t, \mathbf{u}) - g_c(\mathbf{u}).$$

In the continuous time limit, the higher order terms  $O(\Delta t, \mathbf{x}_t, \mathbf{u})$  disappear as these depend on  $\Delta t$ , i.e., i.e.,  $\lim_{\Delta t \rightarrow 0} O(\Delta t, \mathbf{x}_t, \mathbf{u}) = 0$ . The action is also independent of the discounting as  $\lim_{\Delta t \rightarrow 0} \gamma = 1$ . Therefore, the continuous time optimal action is defined as

$$\mathbf{u}_t^* = \arg \max_{\mathbf{u}} \nabla_x V^T \mathbf{B}(\mathbf{x}) \mathbf{u}_t - g_c(\mathbf{u}).$$

This optimization can be solved analytically as  $g_c$  is strictly convex and hence  $\nabla g_c(\mathbf{u}) = \mathbf{w}$  is invertible, i.e.,  $\mathbf{u} = [\nabla g_c]^{-1}(\mathbf{w}) := \nabla \tilde{g}_c(\mathbf{w})$  with the convex conjugate  $\tilde{g}$ . The optimal action is described by

$$\mathbf{B}(\mathbf{x})^T \nabla_x V^* - \nabla g_c(\mathbf{u}) := 0 \Rightarrow \mathbf{u}^* = \nabla \tilde{g}_c \left( \mathbf{B}(\mathbf{x})^T \nabla_x V^k \right).$$

Therefore, the value function update can be rewritten by substituting the optimal action, i.e.,

$$\begin{aligned} V_{\text{tar}}(\mathbf{x}_t) &= r \left( \mathbf{x}_t, \nabla \tilde{g}_c \left( \mathbf{B}(\mathbf{x}_t)^T \nabla_x V^k \right) \right) + \gamma V^k(\mathbf{x}_{t+1}; \psi_k) \\ \text{with } \mathbf{x}_{t+1} &= f \left( \mathbf{x}_t, \nabla \tilde{g}_c \left( \mathbf{B}(\mathbf{x}_t)^T \nabla_x V^k \right) \right). \end{aligned}$$

□

### Experimental Setup

**Systems** The performance of the algorithms is evaluated using the *swing-up* the torque-limited pendulum, cartpole and Furuta pendulum. The physical cartpole (Figure ??) and Furuta pendulum (Figure ??) are manufactured by Quanser (2018). For simulation, we use the equations of motion and physical parameters of the supplier. Both systems have very different characteristics. The Furuta pendulum consists of a small and light pendulum (24g, 12.9cm) with a strong direct-drive motor. Even minor differences in the action cause large changes in acceleration due to the large amplification of the mass-matrix inverse. Therefore, the main source of uncertainty for this system is the uncertainty of the model parameters. The cartpole has a longer and heavier pendulum (127g, 33.6cm). The cart is actuated by a geared cogwheel drive. Due to the larger masses the cartpole is not so sensitive to the model parameters. The main source of uncertainty for this system is the friction and the backlash of the linear actuator. The systems are simulated and observed with 500Hz. The control frequency varies between algorithm and is treated as hyperparameter.

**Baselines** The control performance is compared to Deep Deterministic Policy Gradients (DDPG) (Lillicrap et al., 2015), Soft Actor Critic (SAC) (Haarnoja et al., 2018) and Proximal Policy Optimization (PPO) (Schulman et al., 2017). The baselines are augmented with uniform domain randomization (UDR) (Muratore et al., 2018). For the experiments the open-source implementations of MushroomRL (D'Eramo et al., 2020) are used. We compare two different initial state distributions ( $\mu$ ). First, the initial pendulum angle  $\theta_0$  is sampled uniformly, i.e.  $\theta_0 \sim \mathcal{U}(-\pi, +\pi)$ . Second, the initial angle is sampled from a Gaussian distribution with the pendulum facing downwards, i.e.,  $\theta_0 \sim \mathcal{N}(\pm\pi, \sigma)$ . The uniform sampling avoids the exploration problem and generates a larger state distribution of the optimal policy.

**Reward Function** The desired state for all tasks is the upward pointing pendulum at  $\mathbf{x}_{\text{des}} = \mathbf{0}$ . The state reward is described by  $q_c(\mathbf{x}) = -(\mathbf{z} - \mathbf{z}_{\text{des}})^T \mathbf{Q}(\mathbf{z} - \mathbf{z}_{\text{des}})$  with the positive definite matrix  $\mathbf{Q}$  and the transformed state  $\mathbf{z}$ . For continuous joints the joint state is transformed to  $z_i = \pi^2 \sin(x_i)$ . The action cost is described by  $g_c(\mathbf{u}) = -2\beta \mathbf{u}_{\text{max}} / \pi \log \cos(\pi \mathbf{u} / (2 \mathbf{u}_{\text{max}}))$  with the actuation limit  $\mathbf{u}_{\text{max}}$  and the positive constant  $\beta$ . This barrier

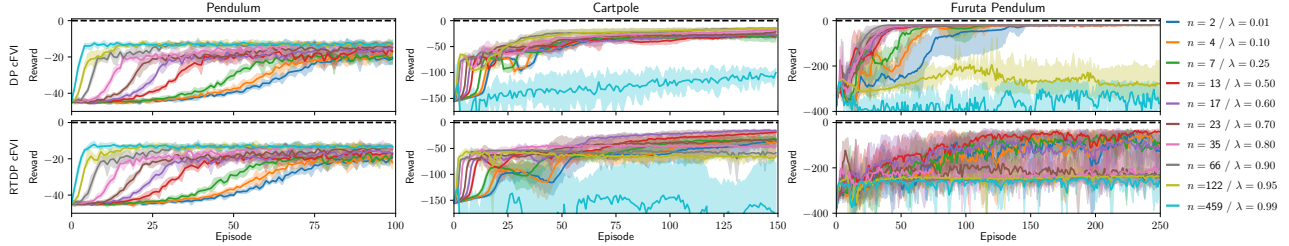


Figure 8. The learning curves averaged over 5 seeds for the  $n$ -step value function target. The shaded area displays the *min/max* range between seeds. The step count is selected such that  $\lambda^n := 10^{-4}$ . Increasing the horizon of the value function target increases the convergence rate to the optimal value function. For very long horizons the learning diverges as it over fits to the current value function approximation. Furthermore, the performance of the optimal policy also increases with roll out length.

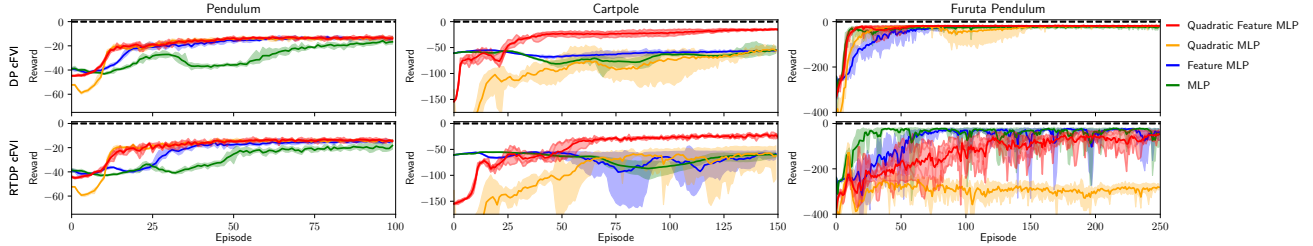


Figure 9. The learning curves averaged over 5 seeds for the different model architectures. The shaded area displays the *min/max* range between seeds. All network architectures are capable of learning the value function and policy for most of the tasks. The locally quadratic network architecture increases learning speed compared to the baselines. The structured architecture acts as an inductive bias that shapes the exploration. The global maximum of the locally quadratic value function is guaranteed at  $\mathbf{x}_{\text{des}}$  and hence the initial policy performs hill-climbing towards this point.

shaped cost bounds the optimal actions. The corresponding policy is shaped by  $\nabla \tilde{g}(\mathbf{w}) = 2 \mathbf{u}_{\text{max}} / \pi \tan^{-1}(\mathbf{w} / \beta)$ . For the experiments, the reward parameters are

$$\text{Pendulum: } \mathbf{Q}_{\text{diag}} = [1.0, 0.1], \quad \beta = 0.5$$

$$\text{Cartpole: } \mathbf{Q}_{\text{diag}} = [25.0, 1.0, 0.5, 0.1], \quad \beta = 0.1$$

$$\text{Furuta Pendulum: } \mathbf{Q}_{\text{diag}} = [1.0, 5.0, 0.1, 0.1], \quad \beta = 0.1$$

**Evaluation** The rewards are evaluated using 100 roll outs in simulation and 15 roll outs on the physical system. If not noted otherwise, each roll out lasts 15s and starts with the pendulum downward. This duration is much longer than the required time to swing up. The pendulum is considered balancing, if the pendulum angle is below  $\pm 5^\circ$  degree for every sample of the last second.

## Extended Experimental Results

### Ablation Study - $N$ -step Value Targets

The learning curves for the ablation study highlighting the importance are shown in Figure 8. This figure contains in contrast to Figure 7 also RTDP cFVI. When increasing  $\lambda$ , which implicitly increases the  $n$ -step horizon (Section 3.3), the convergence rate to the optimal value function increases. This increased learning speed is expected as Equation 7 shows that the convergence rate depends on  $\gamma^n$  with  $\gamma < 1$ .

While the learning speed measured in iterations increases with  $\lambda$ , the computational complexity also increases. Longer horizons require to simulate  $n$  sequential steps increasing the required wall-clock time. Therefore, the computation time increases exponentially with increasing  $\lambda$ . For example the forward roll out in every iteration of the pendulum increases exponentially from 0.4s ( $\lambda = 0.01$ ) to 56.4s ( $\lambda = 0.99$ ). For the Furuta pendulum and the cartpole extremely long horizons of 100+ steps start to diverge as the value function target over fits to the untrained value function. For RTDP cFVI, the horizons must smaller, i.e., 10 - 20 steps. For longer horizons the predicted rollout overfits to the value function outside of the current state distribution, which prevents learning or leads to pre-mature convergence. This is very surprising as even for the true model long time-horizons can be counterproductive due to the local approximation of the value function. Therefore, DP cFVI works best with  $\lambda \in [0.85, 0.95]$  and RTDP cFVI with  $\lambda \in [0.45 - 0.55]$ .

### Ablation Study - Model Architecture

To evaluate the impact of the locally quadratic architecture described by

$$V(\mathbf{x}; \psi) = -(\mathbf{x} - \mathbf{x}_{\text{des}})^T \mathbf{L}(\mathbf{x}; \psi) \mathbf{L}(\mathbf{x}; \psi)^T (\mathbf{x} - \mathbf{x}_{\text{des}}),$$

where  $\mathbf{L}$  is a lower triangular matrix with positive diagonal, we compare this architecture to a standard multi-layer perceptron with and without feature transformation. The learning

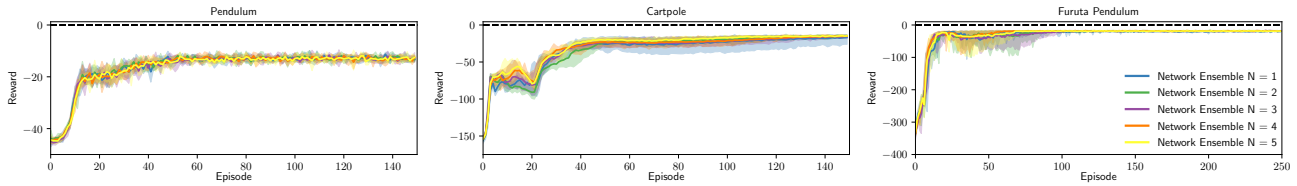


Figure 10. The learning curves averaged over 5 seeds with different model ensemble sizes  $N$ . The shaded area displays the *min/max* range between seeds. The performance of the optimal policy is not significantly affected by the model ensemble. For the cartpole and especially the Furuta pendulum, the larger model ensembles stabilize the training and achieve faster learning and exhibit less variations between seeds.

curves for the ablation study highlighting the importance of the network architecture are shown in Figure 9. The reward curves are averaged over 5 seeds and visualize the maximum range between seeds. For most systems all network architectures are able to learn a good policy. The locally quadratic value function is on average the best performing architecture. The structure acts as an inductive bias that shapes the exploration and leads to faster learning. The global maximum is guaranteed to be at  $x_{\text{des}}$  as  $L(x; \psi)L(x; \psi)^T$  is positive definite. Therefore, the initial policy directly performs hill-climbing towards the balancing position. Only for the cartpole the other network architectures fail. For this system, these architectures learn a local optimal solution of balancing the pendulum downwards. This local optima is the conservative solution as the cost associated with the cart position is comparatively high to avoid the cart limits on the physical system. Therefore, stabilizing the pendulum downwards is better compared to swinging the pendulum up and failing at the balancing. The locally quadratic network with the feature transform, learns the optimal policy for the cartpole. This architecture avoids the local solution as the network structure guides the exploration to be optimistic and the feature transform simplifies the value function learning to learn a successful balancing.

### Ablation Study - Model Ensemble

The learning curves for different model ensemble sizes are shown in Figure 10. The model ensemble does not significantly affect the performance of the final policy but reduces the variance in learning speed between seeds. The variance between seeds also increases. The reduced variance for the model ensembles is caused by the smoothing of the network initialization. The mean across different initial weights lets the initialization be more conservative compared to a single network. For the comparatively small value function networks (i.e., 2-3 layers deep and 64 - 128 units wide), we prefer the network ensembles as the computation time does not increase when increasing the ensemble size. If the individual networks are batched and evaluated on the GPU the computation time does not increase. The network ensembles could also be evaluated at 500Hz for the real-time control experiments using an Intel i7 9900k.