

# LVIS: Learning from Value Function Intervals for Contact-Aware Robot Controllers

Robin Deits<sup>1</sup>, Twan Koolen<sup>1</sup>, and Russ Tedrake<sup>1</sup>

**Abstract**—Guided policy search is a popular approach for training controllers for high-dimensional systems, but it has a number of pitfalls. Non-convex trajectory optimization has local minima, and non-uniqueness in the optimal policy itself can mean that independently-optimized samples do not describe a coherent policy from which to train. We introduce LVIS, which circumvents the issue of local minima through global mixed-integer optimization and the issue of non-uniqueness through learning the optimal value function (or cost-to-go) rather than the optimal policy. To avoid the expense of solving the mixed-integer programs to full global optimality, we instead solve them only partially, extracting intervals containing the true cost-to-go from early termination of the branch-and-bound algorithm. These interval samples are used to weakly supervise the training of a neural net which approximates the true cost-to-go. Online, we use that learned cost-to-go as the terminal cost of a one-step model-predictive controller, which we solve via a small mixed-integer optimization. We demonstrate the LVIS approach on a cart-pole system with walls and a planar humanoid robot model and show that it can be applied to a fundamentally hard problem in feedback control—control through contact.

## I. INTRODUCTION

While there are a variety of successful approaches for planning multi-contact behaviors (e.g. [1], [2], [3], [4]), it has proven to be difficult to apply these techniques quickly enough to be used in response to disturbances. Furthermore, most multi-contact trajectory optimizations are solved via non-convex optimizations, typically through sequential quadratic programming (e.g. [2], [3]) or differential dynamic programming (e.g. [4], [5], [6], [7]). These techniques can generally find locally optimal solutions, but make no guarantees of global optimality. While locally optimal solutions are often sufficient for planning purposes, they make training a policy from examples (as in [7] and [5]) more difficult, as the locally optimal samples may not describe a coherent global policy.

Mixed-integer optimization offers some hope for planning globally optimal multi-contact behaviors: By explicitly representing the discrete changes in dynamics with discrete (i.e. integer) variables, we can create optimization problems which are solvable to global optimality using branch-and-bound [8]. Global optimality is possible even in the presence of nonlinear constraints [9], [10], but for this paper we restrict ourselves to piecewise affine models, inspired by the long history of successful linearized dynamical models for humanoid robots (e.g. [11]). Unfortunately, global optimality

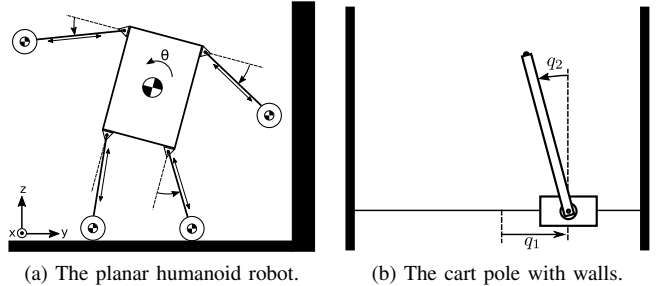


Fig. 1. The robot models used in this work.

comes at a cost, with typical trajectory optimizations taking seconds or minutes to solve [1]. Furthermore, there is no guarantee that these expensive optimizations will result in a consistent global policy, as the optimal policy itself may not be unique.

On the other hand, we do not necessarily need to completely solve a mixed-integer optimization to get some useful information from it. Mixed-integer convex problems are generally solved by branch-and-bound [8], a process which iteratively finds better candidate solutions and tighter bounds on the best possible solution. If we ensure that a candidate solution always exists, then we can terminate the branch-and-bound process at any time, retrieving the best solution and tightest bound found so far. Although we could attempt to train from these sub-optimal solutions, we would again be learning to imitate a sub-optimal controller. The *bounds themselves*, however, are extremely useful: In an MPC problem, bounds on the optimal objective value are also bounds on the optimal cost-to-go<sup>1</sup> from a given state.<sup>2</sup> Having a model of the cost-to-go in turn enables fast online control by simply greedily descending that cost.

## A. LVIS: Learning from Value Interval Sampling

In this work, we introduce LVIS, a new approach for the creation of contact-aware controllers. We model our robot's contact dynamics with complementarity constraints (Sect. IV-A). Offline, we set up a large number of trajectory optimizations in the form of mixed-integer quadratic programs (MIQPs) from a variety of initial states. We partially solve those optimizations, terminating early and extracting

<sup>1</sup>CSAIL, Massachusetts Institute of Technology, Cambridge, MA, USA {rdeits, tkoolen, russt}@csail.mit.edu

This work was supported by the Fannie and John Hertz Foundation and the MIT Computer Science and Artificial Intelligence Lab. Source code is available at <https://github.com/rdeits/LVIS-dev>

<sup>1</sup>The cost-to-go, which we will also refer to as the value function  $J(x)$ , is the cost which will be accumulated by the optimal controller starting from state  $x$  [12].

<sup>2</sup>This assumes that the MPC horizon is long enough to reach a set of states with known cost-to-go. We will violate that assumption later, but attempt to demonstrate empirically that the objective bounds are still useful.

concrete intervals containing the optimal cost at the given robot state (Sect. IV-B). From these intervals, we train a small neural net to approximate the cost-to-go using a loss function which penalizes predicted values outside the known intervals (Sect. IV-C). Online, we run a simple one-step MPC controller to greedily descend the approximate cost-to-go as quickly as possible subject to the robot’s dynamics (Sect. IV-D).

## II. RELATED WORK

This work is similar to that of Zhong et al. in [5], in which offline optimizations were also used to train an approximate cost-to-go used as the terminal cost of a shorter-horizon MPC problem. Zhong’s work differs from ours in its use of iterative LQR (iLQR) to generate the cost-to-go samples. As iLQR is a local nonlinear optimization, it can only provide an estimate of the upper bound of the the cost-to-go (since a lower cost might exist in a space that was not explored by the local optimization). In our case, by constructing a mixed-integer optimization and solving it with branch-and-bound, we recover global upper and lower bounds on cost-to-go, using the interval spanned by those bounds during training. In Sect. V-B.3, we specifically compare LVIS with an approach of learning only from upper bounds on the cost-to-go.

A major obstacle to solving MPC problems for system with contacts is the potentially vast number of possible mode sequences.<sup>3</sup> If an optimal mode sequence for a given state could be computed, then we could perform a cheap continuous optimization to choose the precise optimal input given that mode sequence. This is the approach taken by Hogan in [13], in which a neural net is trained to predict mode sequences from robot states. Marcucci takes a similar approach in [14] by creating a library of provably feasible stabilizing mode sequences and looking up a mode sequence for the robot’s current state at run-time.

Alternatively, the efficient sequential linear-quadratic methods from [6] do allow for locally optimal real-time MPC for systems with contact, avoiding the need for offline learning. These optimizations are still subject to local minima, but the ability to run them at real-time means that they do not need to be used to train a global policy.

Looking more broadly, reinforcement learning offers an alternative approach which does not require any explicit offline planning but instead simply the ability to roll out actions in simulation or hardware (e.g. [15], [16]). In our approach, directly measuring the cost-to-go intervals from a given state, rather than trying to estimate a reward based on expected future actions, allows us to use a very simple supervised learning architecture instead.

## III. ROBOT MODELS

We demonstrate the LVIS controller on two models: a cart-pole system balancing between two walls and a simplified humanoid model. The cart-pole (Figure 1b) consists of an

actuated cart which can accelerate in one dimension, and an unactuated pole which rotates freely. We modify the system by adding two walls on either side (contact is considered only between the tip of the pole and the walls). The cart-pole system has 4 continuous states, 1 input, and 7 discrete modes (contact-free and sticking, sliding up, and sliding down for each wall).

The simplified planar humanoid model (Figure 1a) has 11 degrees of freedom (3 DoF for the planar translation and rotation of the central body, and 2 DoF for the rotation and extension of each limb in the plane). We model contact between each limb and the fixed floor or wall, and we add hard position limits for each of the 8 joints connecting the limbs to the body. The humanoid has 22 continuous states and 11 inputs. Since each of the 4 limbs has 3 states (free, sticking, or sliding in one of 2 directions), and each of the 8 joints has 3 states (free, at its upper limit, or at its lower limit), the system has a total of  $4^4 \times 3^8 = 1,679,616$  discrete modes. We model these modes implicitly using complementarity conditions (Sect. IV-A).

## IV. TECHNICAL APPROACH

### A. Modeling

Following the formulation of Stewart and Trinkle in [17], we model the dynamics of our robot in a contact-implicit manner with complementarity conditions. In discrete time, these dynamics take the form:

$$\mathbf{M}(\mathbf{v}^{l+1} - \mathbf{v}^l) = h\mathbf{f}_{ext}^l + h\mathbf{C} + h\mathbf{B}\mathbf{u}^l \quad (1a)$$

$$\mathbf{q}^{l+1} - \mathbf{q}^l = h\mathbf{v}^{l+1} \quad (1b)$$

where  $h$  is the time step,  $\mathbf{q}^l$  and  $\mathbf{v}^l$  are the system’s generalized configuration and velocity at time step  $l$ ,  $\mathbf{M}$  is the mass matrix,  $h\mathbf{f}_{ext}$  is the external impulse due to contact and friction,  $h\mathbf{C}$  is the impulse caused by Coriolis, and gravitational forces, and  $h\mathbf{B}\mathbf{u}^l$  is the impulse caused by the generalized inputs  $\mathbf{u}$ . Complementarity conditions ensure that there is no force acting at a distance and no sliding frictional force without an accompanying velocity. For example, we constrain the normal force of a contact with a condition of the form:

$$f_{\perp} \perp \phi(\mathbf{q}) \quad (2)$$

where  $f_{\perp}$  is the normal force and  $\phi(\mathbf{q})$  is the separation between the associated contact point and the world. We use the notation  $\mathbf{a} \perp \mathbf{b}$  to mean  $\mathbf{a} \geq 0$  and  $\mathbf{b} \geq 0$  and  $\mathbf{a}^T \mathbf{b} = 0$ , so condition (2) has the effect of ensuring that normal force is nonnegative (no pulling on the ground), separation is nonnegative (no penetration), and contact force can only occur when separation is zero. For the complete set of complementarity conditions used in this work, see [17].

It is important to note that  $\mathbf{M}$ ,  $\mathbf{C}$ ,  $\mathbf{B}$ , and  $\phi$  all depend on the robot’s current state (in general in a nonlinear way). When we write down a trajectory optimization as an MIQP, we cannot represent these nonlinear dependencies, so we currently linearize the dynamics around the current state. This limits our current implementation, as our results are

<sup>3</sup>Our humanoid system has 1,679,616 modes (Sect. III), so a trajectory optimization with a horizon of 10 steps has  $1679616^{10} \approx 1.8 \times 10^{62}$  possible mode sequences, most of which are infeasible.

only valid for the particular linearization. While linearized centroidal dynamics have been dramatically successful in humanoid robotics (e.g. [11]), we hope to explore using the full nonlinear dynamics in future work (see Sect. VI).

### B. Data Collection via Optimal Control

To generate a single sample, we start from some initial state  $\mathbf{x}^0 = [\mathbf{q}^{0\top} \ \mathbf{v}^{0\top}]^\top$ . Given a convex quadratic cost described by matrices  $\mathbf{Q}$ ,  $\mathbf{R}$ , and  $\mathbf{S}$ , we write down a trajectory optimization:

$$\begin{aligned} J^* = & \underset{\substack{\mathbf{x}^1 \dots \mathbf{x}^N, \mathbf{u}^1 \dots \mathbf{u}^N, \\ \mathbf{f}^1 \dots \mathbf{f}^N}}{\text{minimize}} \sum_{l=1}^N [\mathbf{x}^{l\top} \mathbf{Q} \mathbf{x}^l + \mathbf{u}^{l\top} \mathbf{R} \mathbf{u}^l] + \mathbf{x}^{N\top} \mathbf{S} \mathbf{x}^N \\ & \text{subject to} \quad \text{linearized dynamics (1a, 1b)} \\ & \quad \quad \quad \text{complementarity (2)}. \end{aligned} \quad (3)$$

For each scalar complementarity condition  $x_i \perp y_i$  we introduce a new binary variable  $z_i$  and constrain:

$$x_i \geq 0, \quad y_i \geq 0 \quad (4a)$$

$$z_i = 1 \implies x_i = 0 \quad (4b)$$

$$z_i = 0 \implies y_i = 0. \quad (4c)$$

We formulate the implication constraints in (4b, 4c) as linear constraints using a standard big-M formulation [18]. The introduction of the binary variables  $z_i$  converts our optimization into an MIQP, i.e. a program with a quadratic objective, linear constraints, and some variables constrained to take integer values in  $\{0, 1\}$ .

1) *Generating Feasible Solutions as Warm-Starts:* Because the only constraints present in the optimization problem in (3) are those that encode the dynamics and complementarity conditions, any desired behavior of the system must be expressed in the cost matrices  $\mathbf{Q}$ ,  $\mathbf{R}$ , and  $\mathbf{S}$ . This restricts the expressiveness of our optimization, but it also means we can generate *feasible* solutions to the optimization (3) simply by *simulating* the system under any controller subject to the same physical constraints. These simulated trajectories are used as warm-starts in the mixed-integer optimization.

The quality of the feasible solutions we generate depends entirely on the choice of controller used during the simulation. Fortunately, since simulation is computationally cheap compared to a full MIQP solution, we can afford to warm-start with multiple controllers and pick whichever simulation result happens to have lower cost. In our case, we warm-start every optimization by simulating with both a basic LQR controller and the LVIS controller being trained. As we train the LVIS cost-to-go, its performance improves and it tends to provide increasingly good warm-starts.

2) *Early Termination:* Solving the MIQP optimization (3) can be extremely expensive for a robot with as many states and modes as our planar humanoid: a trajectory with just  $N = 10$  steps can take thousands of seconds to solve to

near optimality<sup>4</sup>. The key insight of LVIS is that we do not need to solve all the way to optimality: We can easily generate feasible solutions by simulation, so the process of solving the optimization problem in (3) is a matter of running the branch-and-bound algorithm to iteratively find better solutions and tighter bounds on the optimal cost. At any point, we can simply terminate the optimization and extract the best solution and tightest bound found so far. We label the cost of the best solution found so far (which is an upper bound on  $J^*$  as  $J_{ub}$ , and we label the tightest lower bound on the optimal cost as  $J_{lb}$ .

We somewhat arbitrarily choose to terminate the optimization at a fixed time limit of 3 seconds to balance the suboptimality of each sample with the rate at which samples can be generated (limits of 5 or 10 seconds provided similar performance). Other time limits or termination strategies are certainly possible but have not yet been explored.

### C. Training the Neural Net

The neural net which will approximate our cost-to-go consists of a simple fully-connected feed-forward network with exponential linear unit (ELU) activations [19]. For the humanoid in Fig. 1a we use two hidden layers with 48 units each, while for the cart-pole in Fig. 1b we use two hidden layers with 24 units each. The neural net has a number of input dimensions equal to the number of states in the robot and an output dimension of one. We will label the predicted cost-to-go at a state  $\mathbf{x}$  as  $\hat{J}(\mathbf{x}; \theta)$ , where  $\theta$  are the trainable parameters of the net.

1) *Loss Function:* We train the neural net from a set of training samples, where each sample consists of a tuple of (robot state  $\mathbf{x}$ , cost-to-go lower bound  $J_{lb}$ , and cost-to-go upper bound  $J_{ub}$ ). We penalize the net for predicting values outside of the range  $[J_{lb}, J_{ub}]$  using a double-sided hinge loss, defined as:

$$h(\mathbf{x}, J_{lb}, J_{ub}; \theta) = \begin{cases} J_{lb} - \hat{J}(\mathbf{x}; \theta) & \text{if } \hat{J}(\mathbf{x}; \theta) < J_{lb} \\ 0 & \text{if } J_{lb} \leq \hat{J}(\mathbf{x}; \theta) \leq J_{ub} \\ \hat{J}(\mathbf{x}; \theta) - J_{ub} & \text{if } \hat{J}(\mathbf{x}; \theta) > J_{ub} \end{cases} \quad (5)$$

The total loss is simply the sum of the hinge losses over every sample  $(\mathbf{x}, J_{lb}, J_{ub})$ .

The use of the hinge loss provides a unique advantage: as training proceeds, we may revisit a particular state  $\mathbf{x}$ , and due to a better warm-start we may come up with tighter bounds  $J_{lb}$  and  $J_{ub}$  than we previously discovered. The hinge loss ensures that, so long as the net predicts a value  $\hat{J}$  which falls within the new, tighter bounds, the old, looser bounds do not influence the total loss. If, instead, we were attempting to learn  $J$  by exactly matching the upper or lower bounds or their precise midpoint, then our older samples would tend to pull the net's output away from the newer samples. Fig. 2

<sup>4</sup>For our humanoid robot model, solving to within 1% of the optimal cost required an average of 1160 seconds per trajectory optimization, with some cases taking several hours each. At a nominal control rate of 100 Hz, solving to full optimality would thus require the data collection to run at less than 0.0008% real-time.

shows an example of a learned value function and the interval samples which were used to train it.

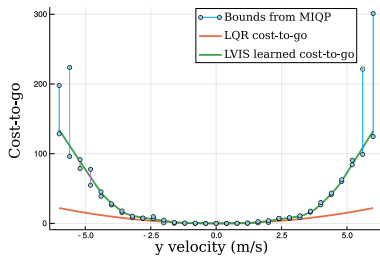


Fig. 2. An example of an approximate value function learned from interval samples, as a function of the planar humanoid’s initial  $y$  velocity. Each connected pair of points represents the interval  $[J_{lb}, J_{ub}]$  from a single trajectory optimization sample. For the states in which the mixed-integer program was solved to optimality,  $J_{lb} \approx J_{ub}$ . The prediction  $\hat{J}(\mathbf{x}; \theta)$  in green lies within the interval  $[J_{lb}, J_{ub}]$  at each sample. The sampled intervals match the LQR cost-to-go (red) for small initial velocities, indicating that the LQR policy is nearly optimal for small disturbances.

2) *Optimization*: The parameters  $\theta$  are trained using a stock ADAM optimizer, with all parameters set to the defaults suggested in [20], and  $\ell^2$  regularization was used to penalize the neural net weights.

#### D. Online Control Using the Learned Cost

The result of the training process described in Sect. IV-C is a neural net whose forward pass approximates the cost-to-go of the original MPC problem. To turn this neural net into a control policy, we construct a new MPC optimization with only one time step ( $N = 1$ ) and set as its terminal cost a local affine approximation of the neural net’s predicted cost-to-go. This corresponds to a greedy gradient descent on the cost-to-go, subject to the robot’s dynamics.

Even the one-step MPC optimization, however, still involves complementarity constraints and is thus a mixed-integer problem. The restriction to a single time step dramatically reduces the number of integer variables which must be solved, allowing near-real-time controller performance, but solving even these smaller MIQPs at control rates is still a challenge. For this work, we implemented a mixed-integer controller in Julia [21] using the `RigidBodyDynamics.jl` software package to model the robot’s dynamics [22], the `Parametron.jl` software package to model the optimization problem [23], and the Gurobi optimizer to solve the resulting problems [24].

#### E. Choosing Initial States with DAGGER

Rather than sampling randomly across the robot’s entire state space, we adopt the DAGGER approach from [25]. Put simply, DAGGER relies on simulating the system using the (initially poorly-trained) policy as its controller instead of the expert, iteratively collecting new training samples from the regions of state-space visited by the learned policy. Our training alternates between (a) letting the learned controller drive the robot for 25-100 time steps while running the mixed-integer optimization to produce new training samples and (b) using those new samples to further train the approximate cost-to-go.

#### F. Policy Net

Rather than trying to learn the value function, we could simply attempt to train a neural net to mimic the mapping from  $\mathbf{x}$  to  $\mathbf{u}$  using the same trajectory optimization samples. We label this approach the Policy Net. As discussed in Sect. IV-B.2, however, the fact that the trajectory optimizations are not generally solved to optimality means that the  $\mathbf{u}$  samples are not generally optimal. Training a neural net to approximate these suboptimal samples is unlikely to result in a good approximation of the optimal policy, but we attempt to do so in order to evaluate that claim.

### V. RESULTS

#### A. Cart-Pole With Walls

The approximate cost-to-go for the cart pole was trained from 3862 mixed-integer trajectory optimization samples. Each trajectory optimization had a horizon of 20 steps and a time step of 25 ms, for a total lookahead time of 0.5 s. Trajectory optimizations were terminated after 3 seconds, which was sufficient for 92.0% of samples to converge to within 1% of the globally optimal cost. As a baseline policy, a discrete-time LQR controller was constructed using the linearization of the system about the upright configuration of the pole. The resulting LQR cost-to-go was used as the terminal cost during the offline mixed-integer trajectory optimization.

Training the cart-pole cost-to-go required approximately two hours on a single CPU, the majority of which was spent solving mixed-integer trajectory optimizations. A total of 500 rounds of training with the ADAM optimizer were performed. Convergence was estimated from 20% of the training samples held as a validation set.

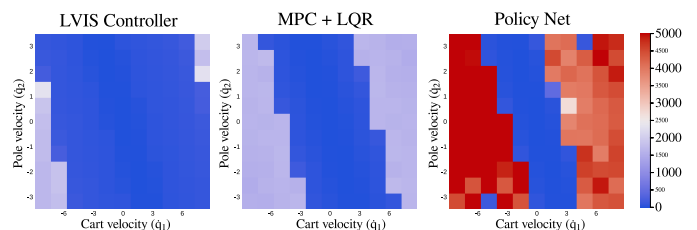


Fig. 3. Accumulated cost of the cart-pole controllers. Each cell indicates the total accumulated cost over a 4-second simulation from the given initial cart and pole velocities, using the same cost matrices as the LQR controller. The regions of very low (dark blue) accumulated cost indicate simulations for which the pole was successfully balanced. The LVIS approach resulted in the lowest accumulated cost and successful stabilization from the widest variety of initial conditions.

1) *Evaluation*: Three potential controllers were evaluated for the cart-pole in order to measure the effectiveness of the learned value function approach:

- 1) LVIS: One-step mixed-integer MPC using the value function learned from the  $[J_{lb}, J_{ub}]$  intervals as its terminal cost.
- 2) MPC + LQR: One-step mixed-integer MPC using the LQR cost-to-go as its terminal cost.

- 3) Policy Net: The neural net trained to mimic the optimal policy (Sect. IV-F).

Each controller was evaluated by simulating the cart-pole for 4 seconds from a range of initial velocities. Each simulation began with the cart centered ( $q_1 = 0$ ) and the pole upright ( $q_2 = 0$ ), with initial cart velocity  $\dot{q}_1$  ranging uniformly from  $-8 \text{ m/s}$  to  $8 \text{ m/s}$  and initial pole rotational velocity  $\dot{q}_2$  ranging uniformly from  $-\pi \text{ rad/s}$  to  $\pi \text{ rad/s}$ . Eleven samples of each initial velocity were collected, for a total of 121 simulations of each controller. Performance of the controller was evaluated by measuring the total accumulated cost (using the same quadratic cost matrices  $Q$  and  $R$  that were used to design the LQR controller) over each simulation.

Results of the cart-pole simulation are shown in Fig. 3. The LVIS controller out-performed the MPC + LQR baseline, resulting in a lower accumulated cost and successful stabilization of the pole from a wider range of initial velocities. The policy net controller (see Sect. IV-F) was able to stabilize the pole from a few initial velocities, but it accumulated more cost than either of the MPC approaches in nearly every case.

### B. Planar Humanoid

For the planar humanoid, approximately 33,700 trajectory optimization samples were collected, with each trajectory optimization having a horizon of 10 and a time step of 50 ms, for a total lookahead of 0.5 s. Trajectory optimizations were terminated after 3 s of optimization with Gurobi. The higher state dimension and larger number of discrete modes made the humanoid trajectory optimization problems substantially harder to solve within that time limit, and only 3.2% of optimizations could be solved to within 1% of the optimal cost within that time.

The method of Mason et al. [26] was used to generate an LQR policy consistent with the contact dynamics of the humanoid (the similar method of [27] could also be used). The LQR policy was designed for the nominal configuration of the robot, shown in the left-most column of Fig. 4, with both feet in contact with the ground.

Training the humanoid value function required approximately 36 CPU hours, again with the majority spent collecting trajectory optimization samples. A total of 300 rounds of training with the ADAM optimizer were performed, and convergence was estimated from 20% of the training samples held as a validation set.

A policy net was also trained on the humanoid optimization samples in an attempt to directly learn the mapping from state to action. The policy net also had two hidden layers with 48 units each, but had 11 outputs, corresponding to the 11 input dimensions of the robot. The same DAGGER training process was run for the policy net, and the same 33,700 samples were provided for training.

1) *Evaluation*: The learned controller was evaluated by simulating the humanoid robot from a variety of initial velocities. From the nominal configuration, the robot's initial linear velocity (along the  $y$  axis of Fig. 1a) was varied from  $-1.5 \text{ m/s}$  to  $1.5 \text{ m/s}$  and its initial angular velocity (about

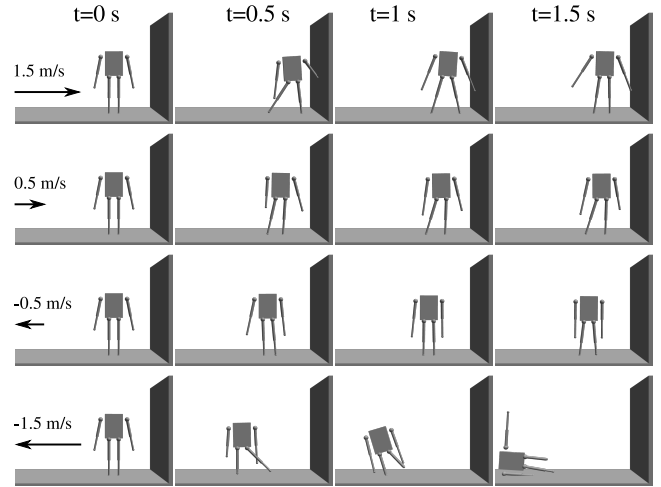


Fig. 4. Animations of the planar humanoid recovering from pushes using the LVIS controller. Initial velocities refer to the velocity of the robot's torso along the  $y$  axis of Fig. 1a. Note that the robot can recover from a  $1.5 \text{ m/s}$  velocity to the right by using contact with the wall, but cannot recover from the same initial velocity to the left.

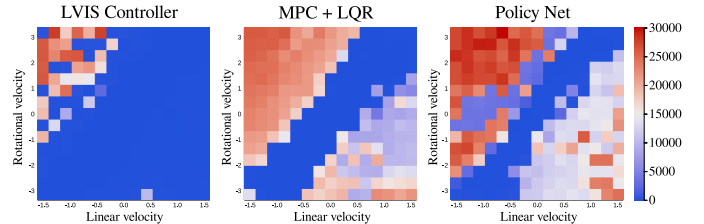


Fig. 5. Accumulated cost of the humanoid controllers. Each cell indicates the total accumulated cost over a 4-second simulation from the given initial linear and angular velocity of the robot's body, using the same cost matrices as the LQR controller. LVIS achieved a low accumulated cost (dark blue) across a wide variety of initial conditions, performing particularly well in the bottom-right corner of the grid in which the initial velocity moved the robot towards the wall.

the  $x$  axis of Fig. 1a) was varied from  $-\pi \text{ rad/s}$  to  $\pi \text{ rad/s}$ . The robot was then simulated under each control policy for 4 seconds using a simulated control rate of 100 Hz.

As was the case with the cart-pole, the LVIS controller generated substantially lower accumulated cost than the baseline controller using the LQR cost-to-go. In particular, the LVIS controller performed especially well when the robot's initial velocity directed it towards the wall, since it was able to both step and reach for the wall in order to maintain balance. The performance of the LVIS controller from a variety of initial velocities can be seen in Fig. 4, and the learned cost is compared with LQR and the Policy Net in Fig. 5.

2) *Capturability Analysis*: Fig. 5 shows that the controller using the learned cost-to-go out-performs the baseline LQR controller, but it could simply be the case that the baseline LQR controller performed very poorly, making it easy to beat. To evaluate the performance of both controllers with respect to an independent benchmark, we can apply the Instantaneous Capture Point (ICP) work of Pratt et al. [28] to estimate the range of initial velocities for which the



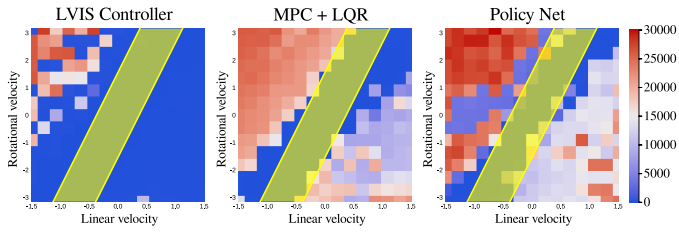


Fig. 6. Comparing the performance of the humanoid controllers with the zero-step capturability region predicted by [28]. Each cell indicates the accumulated cost, as in Fig. 5. For initial velocities in the yellow shaded region, the estimated Instantaneous Capture Point (ICP) lies between the robot’s feet, so it should be possible for a controller to stabilize the center of mass without taking a step. The set of states stabilized by the LQR controller, indicated by the very low (dark blue) accumulated cost, approximately matches the region predicted by the ICP, while the LVIS controller stabilizes a much larger region.

controller should be able to recover without taking a step. Fig. 6 shows in yellow the set of states for which the ICP predicts that the robot should be able to balance without taking a step. Even the baseline LQR controller is able to stabilize the robot from that entire range of states<sup>5</sup>, while the LVIS controller substantially out-performs the capture point predictions, since it is able to exploit changes in contact to stabilize the robot.

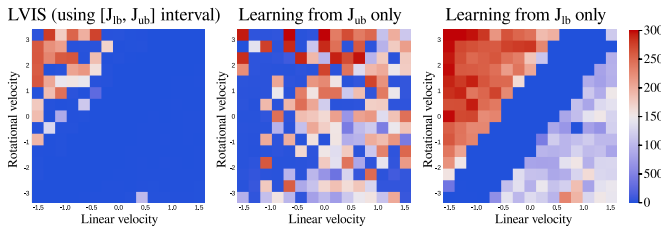


Fig. 7. Comparing learning the cost-to-go from the bounded intervals (Sect. IV-C.1), just the upper bound samples  $J_{ub}$ , or just the lower bound samples  $J_{lb}$ . Accumulated cost was measured as in Fig. 5. Neither the upper nor the lower bounds alone were sufficient to train an approximate cost-to-go which out-performed the LQR baseline or the LVIS approach (left).

3) *The Importance of Intervals:* As described in Sect. IV-C.1, we only penalize the neural net for predicting a cost-to-go which is outside of the interval  $[J_{lb}, J_{ub}]$ . To test the validity of that approach, we also tried training the neural net to exactly mimic the best feasible value  $J_{ub}$  or the best lower bound  $J_{lb}$ . Two additional neural nets were trained using the same neural net structure, number of samples, and training process as in Sect. V-B.1. Each net was penalized for the  $\ell^1$  error between its prediction and  $J_{lb}$  or  $J_{ub}$ , respectively. We evaluated both of these cost-to-go approximations using the same simulation procedure as in Fig. 5. Neither the  $J_{ub}$  samples alone nor the  $J_{lb}$  samples alone produced a cost-to-go and controller which could out-perform even the LQR baseline, as shown in Fig. 7. In practice, attempting to train from just  $J_{lb}$  or  $J_{ub}$  resulted in substantial under-fitting, as the upper and lower bound data both showed a great deal

<sup>5</sup>While the LQR controller has no notion of stepping, it can sometimes successfully slide the foot in the direction of the initial velocity, resulting in stabilization even when the capture point predicts a fall.

of noise from one sample to the next, influenced by the quality of the warm-start solutions, and the varying behavior of Gurobi’s internal heuristics.

### C. Learning in Parameterized Environments

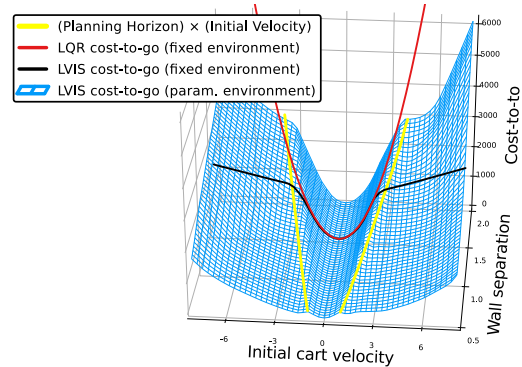


Fig. 8. Learned cost-to-go (blue mesh) as a function of the initial cart velocity  $\dot{q}_1$  and the parameter representing the distance from the center of the track to each wall. The parameterized cost-to-go closely matches LQR (in red) for states for which the cart will not reach the wall within the planning horizon (between the yellow lines), while it rises more gently elsewhere as the walls enable the robot to dissipate energy and reduce its cost. The black line shows the cost function learned in Sect. V-A, for which the wall distance was fixed at 1.5 m.

One drawback of the LVIS approach is that the offline training and trajectory optimization are performed in a fixed environment. This essentially bakes that environment into the learned cost-to-go, resulting in a controller which is only useful in the trained environment. To handle a variety of environments, we suggest creating parameterized templates representing deformable environments. By encoding the environment parameters into the input to the LVIS neural net (both in training and at run-time), we can create a learned cost-to-go which is a function both of the robot’s state and the environment parameters.

As a basic demonstration of this approach, we modified the cart-pole environment (Fig. 1b), adding a single parameter to represent the distance from the center of the track to the walls. The same training process as in Sect. V-A was run, with 20,982 trajectory optimization samples collected over 18 hours. For each iteration of the DAGGER training, the distance to the walls was uniformly randomly varied from 0.5 m to 2.0 m. The resulting learned cost-to-go (now a function of  $x$  and that parameter) is shown in Fig. 8.

## VI. FUTURE WORK

The most significant limitation of this work is the use of the linearized dynamics, which means that the learned cost-to-go is only valid for that linearization. We believe, however, that this limitation can be overcome: there exist spatial branch-and-bound techniques, analogous to the tools used to solve MIQPs, which are applicable to general nonlinear mixed-integer optimizations. These techniques can also provide rigorous upper and lower bounds on the cost-to-go with no need to linearize. We look forward to experimenting with the full nonlinear dynamics, using tools like Couenne [9] or BARON [10] to solve the resulting optimizations.

## REFERENCES

- [1] A. K. Valenzuela, “Mixed-integer convex optimization for planning aggressive motions of legged robots over rough terrain,” Ph.D. dissertation, Massachusetts Institute of Technology, 2016.
- [2] H. Dai, A. Valenzuela, and R. Tedrake, “Whole-body Motion Planning with Simple Dynamics and Full Kinematics,” in *IEEE-RAS International Conference on Humanoid Robots*, Madrid, Spain, 2014.
- [3] M. Posa, C. Cantu, and R. Tedrake, “A direct method for trajectory optimization of rigid bodies through contact,” *The International Journal of Robotics Research*, vol. 33, no. 1, pp. 69–81, Jan. 2014.
- [4] I. Mordatch and E. Todorov, “Combining the benefits of function approximation and trajectory optimization,” in *Robotics: Science and Systems*, 2014.
- [5] 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)*, Apr. 2013, pp. 100–107.
- [6] F. Farshidian, E. Jelavić, A. Satapathy, M. Gifftthaler, and J. Buchli, “Real-Time Motion Planning of Legged Robots: A Model Predictive Control Approach,” *arXiv:1710.04029 [cs]*, Oct. 2017.
- [7] S. Levine and V. Koltun, “Guided Policy Search,” in *ICML*, Atlanta, GA, USA, 2013, p. 10.
- [8] C. A. Floudas, *Nonlinear and Mixed Integer Optimization: Fundamentals and Applications*, ser. Topics in chemical engineering. New York: Oxford Univ. Press, 1995, oCLC: 246760342.
- [9] P. Belotti, J. Lee, L. Liberti, F. Margot, and A. Wächter, “Branching and bounds tightening techniques for non-convex MINLP,” *Optimization Methods and Software*, vol. 24, no. 4-5, pp. 597–634, Oct. 2009.
- [10] M. Tawarmalani and N. V. Sahinidis, “A polyhedral branch-and-cut approach to global optimization,” *Mathematical Programming*, vol. 103, no. 2, pp. 225–249, June 2005.
- [11] S. Kajita, F. Kanehiro, K. Kaneko, K. Fujiwara, K. Harada, K. Yokoi, and H. Hirukawa, “Biped walking pattern generation by using preview control of zero-moment point,” in *IEEE International Conference on Robotics and Automation*, vol. 2, Sept. 2003, pp. 1620–1626.
- [12] R. Tedrake, “Underactuated Robotics: Algorithms for Walking, Running, Swimming, Flying, and Manipulation,” <http://underactuated.csail.mit.edu/underactuated.html>, Mar. 2018.
- [13] F. R. Hogan, E. R. Grau, and A. Rodriguez, “Reactive Planar Manipulation with Convex Hybrid MPC,” in *IEEE International Conference on Robotics and Automation*, May 2018.
- [14] T. Marcucci, R. Deits, M. Gabiccini, A. Bicchi, and R. Tedrake, “Approximate hybrid model predictive control for multi-contact push recovery in complex environments,” in *IEEE-RAS International Conference on Humanoid Robotics*, Birmingham, Nov. 2017.
- [15] J. Kober, J. A. Bagnell, and J. Peters, “Reinforcement Learning in Robotics: A Survey,” *International Journal of Robotics Research*, July 2013.
- [16] S. Gu, E. Holly, T. Lillicrap, and S. Levine, “Deep Reinforcement Learning for Robotic Manipulation with Asynchronous Off-Policy Updates,” *arXiv:1610.00633 [cs]*, Oct. 2016.
- [17] D. Stewart and J. C. Trinkle, “An implicit time-stepping scheme for rigid body dynamics with coulomb friction,” in *Robotics and Automation, 2000. Proceedings. ICRA’00. IEEE International Conference On*, vol. 1. IEEE, 2000, pp. 162–169.
- [18] J. Lofberg, “Big-M and Convex Hulls,” <http://users.isy.liu.se/johanl/yalmip/pmwiki.php?n=Tutorials.Big-MAndConvexHulls>, 2012.
- [19] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, “Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs),” *arXiv:1511.07289 [cs]*, Nov. 2015.
- [20] D. P. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization,” *arXiv:1412.6980 [cs]*, Dec. 2014.
- [21] T. Koolen and R. Deits, “Julia for robotics: Simulation and real-time control in a high-level programming language,” *Submitted to ICRA 2019*, Sept. 2018.
- [22] T. Koolen and contributors, “RigidBodyDynamics.jl,” 2016.
- [23] T. Koolen and R. Deits, “Parametron.jl,” 2018.
- [24] Gurobi Optimization, Inc., “Gurobi Optimizer Reference Manual,” <http://www.gurobi.com/>, 2014.
- [25] S. Ross, G. J. Gordon, and D. Bagnell, “A reduction of imitation learning and structured prediction to no-regret online learning,” in *International Conference on Artificial Intelligence and Statistics*, 2011, pp. 627–635.
- [26] S. Mason, N. Rotella, S. Schaal, and L. Righetti, “Balancing and walking using full dynamics LQR control with contact constraints,” in *IEEE-RAS International Conference on Humanoid Robots*, Nov. 2016.
- [27] M. Posa, S. Kuindersma, and R. Tedrake, “Optimization and stabilization of trajectories for constrained dynamical systems,” in *IEEE International Conference on Robotics and Automation*, May 2016.
- [28] J. Pratt, J. Carff, S. Drakunov, and A. Goswami, “Capture Point: A Step toward Humanoid Push Recovery,” in *IEEE-RAS International Conference on Humanoid Robots*, Dec. 2006.