Neural Lyapunov and Optimal Control

Daniel Layeghi School of Informatics The University of Edinburgh Steve Tonneau School of Informatics The University of Edinburgh Michael Mistry School of Informatics The University of Edinburgh

Abstract—Despite impressive results, reinforcement learning (RL) suffers from slow convergence and requires a large variety of tuning strategies. In this paper, we investigate the ability of RL algorithms on simple continuous control tasks. We show that without reward and environment tuning, RL suffers from poor convergence. In turn, we introduce an optimal control (OC) theoretic learning-based method that can solve the same problems robustly with simple parsimonious costs. We use the Hamilton-Jacobi-Bellman (HJB) and first-order gradients to learn optimal time-varying value functions and therefore, policies. We show the relaxation of our objective results in time-varying Lyapunov functions, further verifying our approach by providing guarantees over a compact set of initial conditions. We compare our method to Soft Actor Critic (SAC) and Proximal Policy Optimisation (PPO). In this comparison, we solve all tasks, we never underperform in task cost and we show that at the point of our convergence, we outperform SAC and PPO in the best case by 4 and 2 orders of magnitude.

I. Introduction

Finding the optimal law to control a non-linear dynamical system with respect to an objective function is an open challenge for many systems. In recent years a class of reinforcement learning (RL) algorithms have empirically demonstrated an ability to to solve complex continuous control tasks [1, 2]. RL does so by leveraging concepts within optimal control (OC) by parameterising and learning the policy or the value function space. Many of these impressive feats however, rely on strategies to improve convergence and stability such as; linearising dynamics regimes via Proportional-Derivative (PD) control or avoiding non-linear rspacesegimes by terminating the episode outside of locally linear state space[3][4]. Additionally, they rely on complex reward-shaping and extensive hyperparameter tuning[5]. Part of the reason RL requires these strategies is that many of the algorithms use zeroth-order gradient estimates to minimise objectives. These estimates are high in variance and can descend into regions of the state space that impair convergence [6].

Trajectory optimisation (TO) is another class of optimal control that has proven to be effective and efficient in problems where dynamics are nonlinear and unstable[4, 7]. The efficiency of TO methods comes from leveraging first or higher-order derivatives of the objective. This results in more stable optimisation that does not require extensive hyperparameters, complex cost design or termination strategies. However, TO methods parametrise trajectory spaces and require reoptimisation per new initial condition.

In this paper, we investigate the possibility of leveraging OCtheoretic tools to learn an optimal policy that allows us to not

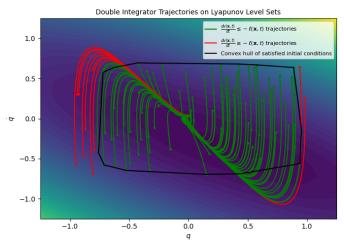


Fig. 1: Compact stability region for double integrator, computed by Neural Lyapunov Control.

require online reoptimisation while avoiding the drawbacks of reinforcement learning. This requires finding a way to leverage Optimal Control theory within massively parallel machine learning frameworks. We ask two questions:

- Can RL algorithms solve simple continuous control tasks without the need for dampening the dynamics, shaping rewards or avoiding regions of the state space?
- Can we solve these problems with minimal hyperparameter optimisation and reward shaping while embracing the entirety of the state and control space?

We believe the above are important problems to consider to improve on the robustness, applicability and reproducibility of learning-based approaches. To this end, we formulate a new OC theoretic function parameterised approach to learn optimal control policies. We utilise differentiable dynamics and the Hamilton-Jacobi-Bellman (HJB) optimality constraint to formulate mathematical programs that enable learning value functions and therefore, optimal policies. Additionally, we show that a specific relaxation of this objective allows us to learn Lyapunov functions, enabling us to further verify the stability of our method over a compact set of initial conditions.

To solve these programs, we leverage neural ODEs [8], a new gradient estimator, and the parallel optimisation capabilities of deep learning. Finally, we compare our method to Soft Actor-Critic (SAC)[9] and Proximal Policy Optimisation (PPO)[10] on selected linear and nonlinear control affine tasks. We employ minimal cost shaping by using simple parsimonious quadratic costs. We do not restrict the landscape of the state

space by early termination. Additionally, we use identical network architectures, episode horizons and boundaries of initial conditions. We empirically show the following:

- RL suffers from poor convergence on environments where minimal reward shaping and environment tuning are used.
- In our experiments our proposed method solves the tasks with significantly faster convergence and variance in random seeds. Outperforming SAC and PPO by at least a factor of 74 and 2 respectively.

II. RELATED WORK

1) RL robustness: Despite RL's effectiveness, several works have studied the shortcomings of the popular approaches. Authors in [5] performed a comprehensive study on the reproducibility of policy gradient algorithms such as PPO, TRPO and DDPG [10, 11, 12]. Their results showed that these algorithms are very sensitive to a variety of parameters. For example, the choice of random seed may lead to significant outperformance on the same solver. The reward is also crucial to the performance, and authors show simple scaling of the rewards may lead to policies that fail entirely. The choice of network architecture and activation function have also been shown to be consequential. The sample complexity of RL based methods has also been in studied. In [13] authors show that policy gradient methods require a high number of samples to converge even on simple linear LQR problems. They also show that applying LQR to simple learnt models outperform policy gradient based methods by orders of magnitude. The source of this sample complexity is also investigated by [6]. Authors show zeroth-order gradient estimates, the underlying estimator in policy gradient methods, are very high variance and inefficient for a variety of continuous control tasks. Thus this raises an important drawback of RL, namely its sample complexity and sensitivity to hyper parameters. As a result, many approaches have been proposed that aim to mitigate these problems.

2) OC theoretic policy/value learning: Approaches to learning OC policy and value functions have been previously explored. Many seminal works exist in offline settings, where off-the-shelf solvers are used to solve OC problems, and the generated data is then used for training. For example [14] applied this formulation to character control. The authors used second-order nonlinear OC to generate state trajectories. A neural network was then used to train a policy that can generate these trajectories and perform complex dynamic behaviour. However, the implicit consequential constraints within the trajectory optimisation problem, such as inverse dynamics, are not considered within the regression. Additionally, this process is computationally expensive and requires compute clusters for solving multiple OC problems in parallel. Other works within the offline setting focus on the value function space. [15] formulate a supervised approach that uses neural networks to learn to fit approximate cost-to-gos, collected from trajectory optimisation. The authors tackle the complex task of humanoid foothold selection. This approach, inspired by [16], is based on the notion that informed terminal constraints and value

functions can reduce planning horizons. However, due to the task complexity and the difficulty of the data generation process, the authors only consider two initial conditions. Additionally, as they solve a supervised problem, implicit constraints such as Bellman backup are not typically considered in the regression.

Work done by [17] addresses the problem of implicit constraints. They do so by approaching the problem from a value iteration perspective. In this case, the one-step HJB policy is used to rollout and compute cost-to-go data that is then fitted to the value function using Bellman backup loss. Authors of [18] also employ a very similar approach, however, the HJB policy is replaced with a trajectory optimiser. Both authors however, do not consider time parameterisation of the value function, even though the data is collected over a finite time. Overall, offline approaches are computationally intensive and primarily rely on supervised training methods that overlook the implicit dynamics of data being learnt.

Other online and more unified approaches have formulated methods that learn Lyapunov and control barrier functions. These methods train neural networks to inherently satisfy Lyapunov stability by minimising the penalty form of Lyapunov constraints. [19] jointly optimises over Lyapunov and a policy network, whilst [20] additionally optimises over a barrier function. Authors in [21] also use a similar approach whilst introducing a verification procedure for certifying the learnt functions. However, these methods do not apply to finite time horizon problems such as trajectory optimisation, as they do not solve for time-varying Lyapunov functions. Similarly the finite horizon formulation of this approach is also investigated in [22]. However, similar to before this approach requires offline data generation by nominal controllers. Perhaps the closest line of work is done by [23] where the authors exploit the structure of neural ODEs with application to policy learning, given known dynamics. Similarly to [23], [8] explores the constrained version of the policy-based formulation. [24] considers the time parameterised version of the policy, whilst introducing further regularisation using offline data. However, the above methods are formulated in policy space and cannot be applied to value and Lyapunov functions synthesis.

More specific HJB based approaches also exist. Authors in [25] focuses on reachability analysis, wherein they relax the Hamilton-Jacobi (HJ) differential equation and employ neural networks to approximate the corresponding value function, effectively mitigating the curse of dimensionality typically associated with classical grid-based solvers. The results however are not compared to any baselines and it is unclear whether backpropogation through time is feasible for high dimensional systems. In a similar vein, [26] utilises a composite loss involving the HJB equation, Hamiltonian, and trajectory cost. This loss is minimised by learning both value and policy functions. The efficacy of this method is demonstrated through comparisons with model-based RL. To best of our understanding this method both learns a policy and a value function independently so the effect of each network is unclear. In [27] authors also consider a similar domain where they learn the costate vector in the maximum principle

formulation. However, this vector only considers the value function in its time differential space making the learnt function limited to the trajectories on which it was trained.

In this work, we introduce a unified approach that solves finite-time parameterised OC problems while leveraging neural ODEs to learn the corresponding time-varying Lyapunov and value functions. We test our method on linear and nonlinear control-affine systems. We compare to RL baselines outperforming all by orders of magnitude.

III. PRELIMINARIES

A. Optimal Control:

Optimal control is grounded within the Hamilton-Jacobi-Bellman (HJB) equation.

$$-v_{t}(\mathbf{x}(t),t) = \ell(\mathbf{x}(t),\mathbf{u}(t)) + \mathbf{f}(\mathbf{x}(t),\mathbf{u}(t))^{\top}v_{\mathbf{x}}(\mathbf{x}(t),t)$$
(1)

where $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u})$ represents the deterministic dynamics, given state $\mathbf{x} \in \mathbb{R}^n$ and control $\mathbf{u} \in \mathbb{R}^m$. $\ell(\mathbf{x}, \mathbf{u})$ represents the state-control cost, $v_t(\mathbf{x}(t), t)$ and $v_{\mathbf{x}}(\mathbf{x}(t), t)$ represent the Legendre notation for partial derivatives of the value function with respect to time and state, respectively. This is a Partial Differential Equation (PDE) defining the time evolution of the value function. In other words the compact description of the relationship between optimal costs for different states and times. In the case of the existence of a C^1 value function, cost and affine dynamics, one can analytically compute a closed loop feedback law optimal for all initial conditions $\mathbf{x}(0)$. The following shows this optimal control for a controlaffine dynamical system where $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}(t)) + \mathbf{g}(\mathbf{x}(t))\mathbf{u}(t)$. Given a quadratic regularisation of control $\|\mathbf{u}\|_{\mathbf{R}}$. Where \mathbf{R} is a positive definite matrix.

$$\min_{\mathbf{u} \in \mathbf{U}} \left[v_t(\mathbf{x}(t), t) + \ell(\mathbf{x}(t), \mathbf{u}(t)) + \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t))^\top v_{\mathbf{x}}(\mathbf{x}(t), t) \right]$$

$$\pi^*(\mathbf{x}(t), t) = -\frac{1}{2} \mathbf{R}^{-1} \mathbf{g}(\mathbf{x}(t))^\top v_{\mathbf{x}}(\mathbf{x}(t), t)$$

The above policy gives us a time-varying optimal policy $\pi^*(\mathbf{x}(t),t)$ for any initial condition. This is an important property of the HJB equation which relies on the strong assumption of the existence of a known value function. On the other hand, Pontryagin Minimum Principle (PMP) alleviates this problem by operating in the tangential space. It converts this PDE into an Ordinary Differential Equation (ODE) by interpreting $v_{\mathbf{x}}$ the gradient of the value as a stand-alone vector \mathbf{p} known as a costate vector. Referring to the right-hand side of equation 2 as Hamiltonian H, the PMP aims to minimise H over a horizon T using an optimal control trajectory \mathbf{u}^* . This trajectory is determined given an initial condition $\mathbf{x}(0)$ and a terminal condition $\ell_{\mathbf{x}}(\mathbf{x}_T) = \mathbf{p}(T)$, subject to the following:

$$\dot{\mathbf{x}} = H_{\mathbf{p}}|_{*}, \qquad \dot{\mathbf{p}}^{*} = -H_{\mathbf{x}}|_{*}$$

$$\mathbf{u}^{*}(t) = \underset{\mathbf{u} \in \mathbf{U}}{\arg \max} H(\mathbf{x}^{*}(t), \mathbf{u}, \mathbf{p}^{*}(t))$$
(3)

Due to its tractability, this framework underpins the majority of trajectory optimisation algorithms such as iterative LQR [4] or differential dynamic programming (DDP) [28]. However,

it only generates an open loop trajectory valid for a single initial condition. This can be a significant limitation as new initial conditions require reoptimisation. For further discussion around these topics, we refer the reader to [29].

B. Neural ODEs

Neural ODEs are both practically and theoretically[30] fundamental to our work. In this section, we briefly describe neural ODEs from an OC perspective. The general objective for neural ODEs is defined as:

$$L(\mathbf{x}(T)) = L\left(\mathbf{x}(0) + \int_{t_0}^{T} \mathbf{f}(\mathbf{x}(t), \mathbf{u}, t) dt\right)$$
(4)

Neural ODEs differ from traditional OC methods in two ways. Firstly, the control parameter \mathbf{u} is a time-invariant variable interpreted as network weights. Secondly, the loss is evaluated at the end of the trajectory due to the analogous nature of the integration time step to hidden layers in a standard neural network. As a result, the canonical equations of neural ODEs are a special case of PMP, with the costate vector $\dot{\mathbf{p}}^* = -H_{\mathbf{x}}|_* = -\mathbf{f}_{\mathbf{x}}^{\top}(\mathbf{x}(t), \theta, t)\mathbf{p}(t)$ ignoring the effects of running loss $L_{\mathbf{x}}(\mathbf{x}(t))$. Therefore, although neural ODEs create a familiar grounding to OC they cannot be immediately applied to problems where a change of state over time matters. The authors of [23] provide a reformulation of neural ODE's gradient estimator for policy learning, which we show can be extended to the value and Lyapunov space in the next section.

IV. LEARNING LYAPUNOV AND VALUE FUNCTIONS

Our approach provides a straightforward and effective framework for learning value and Lyapunov parameterised OC problems while respecting implicit OC constraints.

A. Value functions

Let us focus on the HJB equation 2 under the optimal policy $\pi^*(\mathbf{x}(t),t)$. By moving the inner product between $v_{\mathbf{x}}$ and the dynamics $\mathbf{f}(\mathbf{x}(t),\mathbf{u}(t))$ or $\frac{d\mathbf{x}(t)}{dt}$ to the left-hand side we can rewrite HJB as a definition for the total rate of change of the value function

$$\frac{\partial v(\mathbf{x}(t), t)}{\partial t} + \frac{\partial v(\mathbf{x}(t), t)}{\partial \mathbf{x}(t)} \frac{d\mathbf{x}(t)}{dt} = -\ell(\mathbf{x}(t), \pi^*(\mathbf{x}(t)))
\frac{dv(\mathbf{x}(t), t)}{dt} = -\ell(\mathbf{x}(t), \pi^*(\mathbf{x}(t)))$$
(5)

The above can be interpreted as a constraint on the rate of change of the value function under the optimal policy. Integrating over a horizon T allows us to evaluate the consistency of optimal policy with respect to the above constraint. This leads to:

$$\int_0^T \frac{dv(\mathbf{x}(t), t)}{dt} = -\int_0^T \ell(\mathbf{x}(t), \pi^*(\mathbf{x}(t)))$$

$$v(\mathbf{x}(T), T) - v(\mathbf{x}(0), 0) = -\int_0^T \ell(\mathbf{x}(t), \pi^*(\mathbf{x}(t)))$$
(6)

Equation 6, provides us with an equality that defines the evolution of the value function over the horizon T as a function

of the running loss. However, we aim to learn the granular temporal and spatial change in the value function over the full horizon. In order to capture this effect we discretise both integrals over Δt increments and apply the left Riemann sum approximation.

$$\int_{t}^{t+\Delta t} \frac{dv(\mathbf{x}(t), t)}{dt} = -\int_{t}^{t+\Delta t} \ell(\mathbf{x}(t), \pi^{*}(\mathbf{x}(t), t))$$
$$v(\mathbf{x}_{t+\Delta t}, t + \Delta t) - v(\mathbf{x}_{t}, t) \simeq -\Delta t \times \ell(\mathbf{x}_{t}, \pi^{*}(\mathbf{x}_{t}, t))$$
(7)

equation 7 provides an approximate temporal and spatial instantaneous constraint on the value function. By rearranging and squaring we can convert this constraint into a soft penalty P_v where:

$$P_{v} = \left(v\left(\mathbf{x}_{t+\Delta t}, t + \Delta t\right) - v\left(\mathbf{x}_{t}, t\right) + \Delta t \times \ell\left(\mathbf{x}_{t}, \pi^{*}\left(\mathbf{x}_{t}, t\right)\right)\right)^{2}$$
(8)

The mathematical program 9 leverages this penalty to learn an approximate value function $\tilde{v}(\mathbf{x},t;\theta)$ that minimises the integration of this penalty over the horizon T and a compact set of initial conditions $\mathbf{x}_0 \in X_0$ where $|X_0| = K$. Program 9 uses the optimal feedback policy shown in 2 and is therefore valid only for control-affine dynamics and quadratic regularisation of controls. However analytical optimal policies may be computed under any convex and PSD control constraints and are therefore easily incorporated within this framework. This is further discussed in Section VI.

$$\min_{\theta} \left[\frac{1}{K} \sum_{k=0}^{K} \sum_{n=0}^{N-1} \left(\tilde{v} \left(\mathbf{x}_{n+1,k}, n+1; \theta \right) \right. \right. \\
\left. - \tilde{v} \left(\mathbf{x}_{n,k}, n; \theta \right) + \Delta t \times \left(\| \mathbf{u}_{n,k}^* \|_{\mathbf{R}} + \ell \left(\mathbf{x}_{n,k} \right) \right) \right)^2 \right]$$
s.t:
$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t)) + \mathbf{g}(\mathbf{x}(t)) \mathbf{u}^*(\mathbf{x}, t) \\
\mathbf{u}^*(\mathbf{x}(t), t) = -\frac{1}{2} \mathbf{R}^{-1} \mathbf{g}(\mathbf{x}(t))^{\top} v_{\mathbf{x}(t)}(\mathbf{x}(t), t; \theta)$$

The above mathematical program defines a finite-time parameterised OC problem that aims to learn a value function that approximately minimises the Bellman backup condition in equation 7 over horizon T, where $N=\frac{T-t_0}{\Delta t}$ represents the number of discrete time steps in the interval from t_0 to T. Each $t_n=n\Delta t+t_0$ and $t_{n+1}=(n+1)\Delta t+t_0$ specify the time instances at the n-th and n+1-th time steps, respectively. Here, $\mathbf{x}_{n,k}$ and $\mathbf{u}_{n,k}^*$ denote the state and control input at time t_n .

B. Lyapunov functions

We briefly motivate constructing Lyapunov programs. In program 9 our objective is to minimise the square Bellman error. While it is possible to minimise, approximation errors may lead to unquantifiable loss of performance in trajectory cost. The Lyapunov function relaxes this constraint to an inequality, thereby, allowing us to convert the problem to a satisfability condition. As a result, we are able to provide performance

guarantees on compact regions of the state space, which are not possible for program 9.

We now focus on the derivation of the Lyapunov program. Finite-time Lyapunov analysis states that if we have a continuously differentiable positive definite function $v(\mathbf{x}(t),t)$, where $v(\mathbf{x}(t),t) > 0$ for $\mathbf{x}(t) \neq 0$ and v(0,t) = 0 then $\mathbf{f}(\mathbf{x}(t),t)$ is stable if: $\dot{v}(\mathbf{x}(t),t) = \frac{\partial v(\mathbf{x}(t),t)}{\partial \mathbf{x}(t)} \mathbf{f}(\mathbf{x}) + \frac{\partial v(\mathbf{x}(t),t)}{\partial t} < 0$, $\dot{v}(0,t) = 0$ 0. This condition must hold for all x and all t. However, finding such Lyapunov functions is not trivial and is unknown apriori. Equation 5 and the conditions above impose constraints on the rate of change of function v in both the contexts of value and Lyapunov functions. However, the Lyapunov condition is a relaxed version of the HJB equality constraint. Constraint 5 is a lower bound where satisfying it transforms v from a Lyapunov function into an optimal value function. We aim to learn an approximate Lyapunov function that allows for generating controls that can sub-optimally complete a task but also satisfy Lyapunov constraints within a compact set of initial conditions. To formulate this we relax the HJB condition and define the Lyapunov constraint with respect to a task loss $\ell(\mathbf{x})$.

$$\frac{\partial v(\mathbf{x}(t), t)}{\partial t} + \frac{\partial v(\mathbf{x}(t), t)}{\partial \mathbf{x}(t)} \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)) \le -\ell(\mathbf{x}(t), \mathbf{u}(t))$$

$$\frac{dv(\mathbf{x}(t), t)}{dt} \le -\ell(\mathbf{x}(t), \mathbf{u}(t))$$

Performing the same integration over discrete time step Δt we obtain an instantaneous inequality constraint that can be converted to a soft penalty P_L :

$$P_{L} = \max \left(v \left(\mathbf{x}_{t+\Delta t}, t + \Delta t \right) - v \left(\mathbf{x}_{t}, t \right) + \Delta t \times \ell \left(\mathbf{x}_{t}, \mathbf{u} \left(\mathbf{x}_{t}, t \right) \right), 0 \right)$$
(11)

The program 12 uses this penalty to formulate an objective function that aims to learn an approximate Control Lyapunov Function (CLF) that enforces this Lyapunov constraint over a horizon T and K set of initial conditions.

$$\min_{\theta} \left[\frac{1}{K} \sum_{k=0}^{K} \sum_{n=0}^{N-1} \max \left(\tilde{v} \left(\mathbf{x}_{n+1,k}, n+1; \theta \right) \right. \right. \\
\left. - \tilde{v} \left(\mathbf{x}_{n,k}, n; \theta \right) + \Delta t \times \left(\| \mathbf{u}_{n,k}^* \|_{\mathbf{R}} + \ell \left(\mathbf{x}_{n,k} \right) \right), 0 \right) \right]$$
s.t:
$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t)) + \mathbf{g}(\mathbf{x}(t)) \mathbf{u}^*(\mathbf{x}, t) \\
\mathbf{u}^*(\mathbf{x}, t) = -\frac{1}{2} \mathbf{R}^{-1} \mathbf{g}(\mathbf{x})^{\top} v_{\mathbf{x}}(\mathbf{x}, t; \theta)$$
(12)

We make use of the same HJB feedback policy used in value learning. This is because, as we show in the Appendix, the policy is a valid CLF for a given Lyapunov function under the same affinity and regularisation assumptions as program 9. We show this in Theorem 1 of the supplementary material. However, other choices for a CLF such as Sontag's universal controller [31] are also valid. Subscripts n and k represent the same parameters as program 9. To parameterise the Lyapunov function while respecting Lyapunov constraints we make use of input convex neural networks (ICNNs) [32] with an additional

positive constant term to approximate the positive definite Lyapunov function: $v(\mathbf{x}, t; \theta) = v_{\text{ICNN}}(\mathbf{x}, t; \theta) + \epsilon ||\mathbf{x}||_2^2$.

1) Neural ODE: The mathematical programs 9 and 12 have a bounded time horizon. While the neural ODE gradient estimator are suitable for finite horizon optimisation, they cannot be directly applied to these programs as the task loss needs to be evaluated over every timestep of the entire horizon T. To alleviate this, we follow the same modifications shown in [23]. These modifications aim to convert the canonical equations of neural ODEs to the PMP case.

$$\dot{\mathbf{a}} = -\mathbf{a}(t)^{\top} \mathbf{f}_{\mathbf{x}(t)}(\mathbf{x}, \mathbf{u}(\mathbf{x}(t), t, \theta)) - \ell_{\mathbf{x}}(\mathbf{x}(t), \mathbf{u}(\mathbf{x}(t), t; \theta)) \dot{\ell}_{\theta} = -\mathbf{a}(t)^{\top} \mathbf{f}_{\mathbf{u}}(\mathbf{x}(t), \mathbf{u}(\mathbf{x}(t), t, \theta)) \mathbf{u}_{\theta}(\mathbf{x}(t), t, \theta) - \ell_{\mathbf{u}}(\mathbf{x}(t), \mathbf{u}(\mathbf{x}(t), t; \theta)) \mathbf{u}_{\theta}(\mathbf{x}(t), t, \theta)$$
(13)

These canonical equations differ from vanilla neural ODEs mentioned in [30]. The resulting gradient estimator is less efficient in complexity compared to its vanilla form. However, it is a necessary addition.

We have shown how parameterised OC problems can be formulated using the mathematical programs 12 and 9. In the next segment, we will demonstrate how these problems can be solved using the gradient estimator 13. We obtain these results on systems with control-affine dynamics.

V. EMPIRICAL RESULTS

1) Criteria: To evaluate our method, we apply it to tasks with linear or locally linear dynamics, such as the Double Integrator and Cartpole stabilisation. We also consider nonlinear problems like Cartpole swing-up and the planar reacher environment. Our method successfully learns time-varying Lyapunov and value functions that approximately satisfy the Bellman and Lyapunov constraints through minimising programs 9 and 12. For a comprehensive assessment, we compare the effectiveness of the learned functions against Soft Actor Critic (SAC) [9] and Proximal Policy Optimisation (PPO) [10] using the Stable-Baselines3 implementation [33]. The hyperparameters for the baseline experiments can be found in our supplementary material. The results are presented in accompanying table I and figure 2. For all tasks, we used our implementation of neural ODEs gradient estimator defined in equation 13. The results were obtained on a core i9 Intel processor with Nvidia GeForce RTX 4070 GPU. The Cartpole and the Reacher models are adopted from [34][35]

2) Environments and solver setting:

Dynamics: Since our method leverages first-order gradients, we require gradients to be available within the PyTorch graph. As a result, we implement our own dynamics.

In the introduction we mentioned convergence strategies such as dampening dynamics and avoiding areas of the state and control space. Rigid body dynamics is defined by $\mathbf{M}(\mathbf{q})\ddot{\mathbf{q}} + \mathbf{C}(\mathbf{q},\dot{\mathbf{q}})\dot{\mathbf{q}} + \mathbf{b}(\dot{\mathbf{q}}) + \mathbf{G}(\mathbf{q}) = \tau$. However, in the case of very high friction: $\mathbf{C}(\mathbf{q},\dot{\mathbf{q}})\dot{\mathbf{q}} \ll \mathbf{b}(\dot{\mathbf{q}})$, and the resulting dynamics becomes $\mathbf{M}(\mathbf{q})\ddot{\mathbf{q}} + \mathbf{b}(\dot{\mathbf{q}}) + \mathbf{G}(\mathbf{q})$. Thus by tuning the environment with high friction coefficients, one can practically

remove a significant nonlinearity, as done in the Mujoco Reacher environment used in OpenAI's Gym [36]. With regards to the state space, for example in CartPole-V1, the control space is discrete where; 0 is push cart to the left and 1 is push cart to the right. Also, the termination condition is active not far from the upright where; $|\mathbf{q}_{pole}| \geq 0.2$ rads and $|\mathbf{q}_{cart}| \geq 2.4$. For the Inverted Double Pendulum; the control space is constrained within $-1 \le \mathbf{u} \ge 1$. The termination condition is active when the pole length falls below 1 and the upright height is 1.196; $l_1 \cos(\mathbf{q}_{pole1}) + l_2 \cos(\mathbf{q}_{pole1} + \mathbf{q}_{pole2}) \leq 1$. Similar termination conditions are also used on the continuous control CartPole-V4 environments. Although such strategies can help convergence, they lead to policies that are trained on small regions of the state space. Additionally, boundaries for such constraints are not clear apriori. For instance control constraints can directly impact the behaviour of the policy and in many cases we do not know the control bounds at which the optimal solution is achieved. Additionally, dampening dynamics can also linearise and simplify the problem. However, it can lead to policies that are trained on unrealistic dynamics. To evaluate the complete effectiveness of the methods we avoid such strategies in our experiments. In the next section we explain the details of our reward/cost design, control and space boundaries and termination conditions.

Solver setting: We aim to keep the setting identical between solvers. As a result, timestep Δt , horizon/episode length T, number of timestep, and rewards/cost are equal between solvers. Additionally, we make the control space unbounded $\mathbf{u} \in \mathbb{R}$, however, we regularise this space with quadratic regularisation $\|\mathbf{u}\|_{\mathbf{R}}$, \mathbf{R} is also identical between solvers per task. Our reward design is the negative of the cost function. Where our cost functions are the simple sum of running quadratic functions, $\cos t = \mathbf{x}^{\top}(t)\mathbf{Q}\mathbf{x}(t) + \mathbf{x}^{\top}(T)\mathbf{Q}_T\mathbf{x}(T)$. We also do not employ any state-based termination conditions and only rely on the horizon to terminate. We do however try to tune the hyperparameters of each baseline solver to the best of our ability, starting from the parameters mentioned in the original papers [10] [9].

A. Value Function Results

1) Reacher: For this fully actuated system, the goal state is defined as $\mathbf{x} = [0, 0, 0, 0]$. It is important to mention the dynamics here follows the low friction definition mentioned in Section V-2. The cost function is quadratic with respect to $\mathbf{x} \in \mathbb{R}^4$ and control $\mathbf{u} \in \mathbb{R}^1$ with $\mathbf{Q} = \text{diag}(1, 1, 0, 0)$ and $\mathbf{R} = \mathbf{M}^{-1}(\mathbf{q})$ where $\mathbf{M}^{-1}(\mathbf{q})$ is the inverse inertial matrix. Additionally we use a terminal cost $\mathbf{Q}_T = \text{diag}(100, 100, 1, 1)$. The negative of this cost is used as a reward. The value function is fully connected neural network (FCN) $V: \mathbb{R}^{4+1} \to \mathbb{R}^1$: (5-128-128-1 FCN). All discretisation timesteps are 0.01s and the Adam optimiser for training [37]. As shown in table I, we satisfy the HJB constraint 5 with an error of 0.25 \pm 0.10. Constraint satisfaction approximately converges at 5e4, however, small errors result in task cost converging at approximately 1e5 timesteps. Our average final cost outperforms SAC and PPO by a factor of 7.43 and 1836.65.

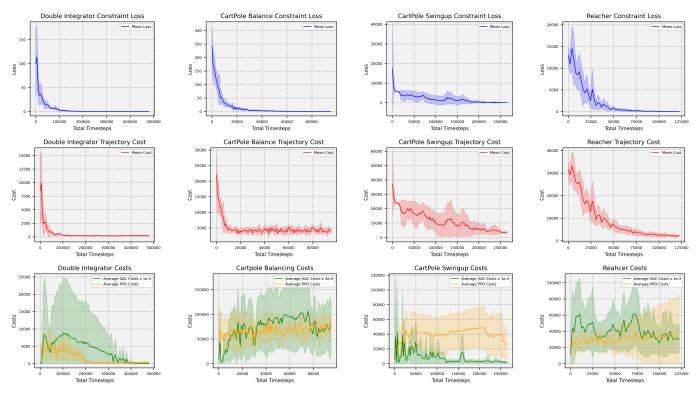


Fig. 2: Top row: Constraint satisfaction loss for value and Lyapunov function constraints. Middle row: Trajectory cost using our method. Bottom row: SAC and PPO trajectory cost. **Due to high values, SAC costs are scaled for visualisation**.

Environment	Constraint Loss	Trajectory Cost			Cost Improvement		Horizon
		Ours	PPO	SAC	PPO	SAC	
Reacher	33.61 ± 20.25	2086.24 ± 419.84	$15514.49, \pm 4465.52$	3831267.56 ± 471973.23	7.43	1836.65	170
Swing up	24.65 ± 16.29	3568.74 ± 887.16	33147.12 ± 16041.53	$255854.68, \pm 81662.83$	144.47	10379.47	171
Balancing	0.07 ± 0.12	3868.80 ± 1012.27	93109.16 ± 15629.57	$2360774.63, \pm 1098978.99$	24.07	284.12	79
Double Integrator	0.03 ± 0.03	178.92 ± 77.63	$359.81, \pm 128.76$	$13311.58, \pm 7681.00$	2.01	74.78	400

TABLE I: Training statistics and performance comparison against SAC and PPO

2) Cartpole Swing Up: The Cartpole provides an underactuated environment obtained from [34]. In this, the goal is defined at the unstable equilibrium $\mathbf{x} = [0,0,0,0]$. A the quadratic cost function is used with respect to $\mathbf{x} \in \mathbb{R}^4$ and control $\mathbf{u} \in \mathbb{R}^1$ with $\mathbf{Q} = \mathrm{diag}(0,0,0,0)$ and $\mathbf{R} = \mathbf{M}^{-1}(\mathbf{q})$. Additionally, we use a terminal cost $\mathbf{Q}_T = \mathrm{diag}(80,600,.8,4.5)$. The negative of this cost was used for baselines. The value function is parameterised by $\tilde{V}_{tl} : \mathbb{R}^{4+1} \to \mathbb{R}^1$: (5-128-128-1 FCN). The program was solved using the same timestep and optimiser as Reacher. The program satisfied HJB constraint by 24.65. Constraint satisfaction approximately converges at 2e5 with small errors causing cost convergence at under 2.5e5 timesteps. As shown in table I SAC and PPO show no convergence in this time.

B. Lyapunov Function Results

1) Double Integrator: We also assess the ability of program 12 to satisfy trajectory stability Lyapunov constraints on locally linear systems. We consider a fully linear double integrator with the goal state at $\mathbf{x} = [0,0]$. For this problem we require the rate of change of the stable trajectories to be upper bounded by the quadratic cost with respect to $\mathbf{x} \in \mathbb{R}^2$ and control

 $\mathbf{u} \in \mathbb{R}^1$ with $\mathbf{Q} = \mathrm{diag}(10,0.1)$, $\mathbf{Q}_T = \mathrm{diag}(10,0.1)$ and $\mathbf{R} = \mathbf{M}^{-1}(\mathbf{q})$ with the Lyapunov function $\tilde{V}: \mathbb{R}^{2+1} \to \mathbb{R}^1$: (3-64-64-1 ICNN). Our results show that we satisfy the Lyapunov constraint 10 with error 0.03 ± 0.03 . Additionally, we show that if we keep the set of initial conditions \mathbf{X}_0 we are able guarantee Lyapunov stability of 90% of the compact set of $\mathbf{K} = 100$ initial conditions. This is also shown in the figure 1. Similar to the previous cases constraint satisfaction loss converges quicker than trajectory cost. Surprisingly SAC significantly underperforms on this task. We outperform SAC and PPO by a factor of 74.78 and 2.01.

2) Cartpole Balancing: We apply program 12 locally linear task, Cartpole balancing. Again we do not terminate episode based on any state V-2. Our termination is only based on the end of Horizon. The goal state is at $\mathbf{x} = [0,0]$ with initial conditions $\mathbf{X}_0 \in (-0.6,0.6)$. Similarly, rate of change of the stable trajectories to be upper bounded by the quadratic cost with respect to $\mathbf{x} \in \mathbb{R}^4$ and control $\mathbf{u} \in \mathbb{R}^1$ with $\mathbf{Q} = \mathrm{diag}(0,25,0.5,0.1)$, $\mathbf{Q}_T = \mathrm{diag}(0,25,0.5,0.1)$ and $\mathbf{R} = \mathbf{M}(\mathbf{q})^{-1}$ with the Lyapunov function $\tilde{V}: \mathbb{R}^{4+1} \to \mathbb{R}^1$: (5-200-500-1 ICNN). In this, the Lyapunov constraint is satisfied

10 with error 0.17 \pm 0.16. Additionally, the same program can satisfy Lyapunov stability of 85% of the compact set of K = 100 constant initial conditions. Constraint and task cost both converge at approximately 3e4 timesteps, with the final cost reaching 3868.80, outperforming SAC and PPO by a factor of 284.12 and 24.07.

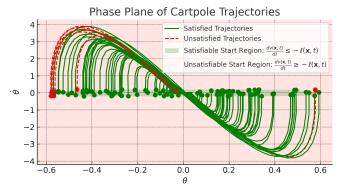


Fig. 3: Cartpole balancing Lyapunov trajectories.

VI. DISCUSSION AND FUTURE WORK

The results show that the proposed programs notably outperform both SAC and PPO in terms of: convergence, costs associated with generated solutions, variance in results, and training stability. This difference is largely attributed to leveraging dynamics and its derivative information within training, analytically encoding dynamic structure within the policy, deterministic dynamics, and time parameterisation of the policy. In our evaluation, we treated the dynamics models as exact. This is a widespread practice within OC, because of the certainty equivalence assumption. This is where state feedback and re-optimisation through methods like Model Predictive Control can compensate for surprising amounts of error [13]. Derivative information through optimisation allows for considering the dynamic effects when minimising task cost. Additionally, the analytical inclusion of model parameters essentially compensates for the dynamic effects of the system, and reduces the complexity of the search space. Typically, RL algorithms are optimised over fixed horizons or episode lengths but do not consider this time dependency within the policy. This essentially results in using an infinite-horizon policy for a finite-time task, which is provably sub-optimal.

There are also challenges with the proposed programs. The HJB constraint and its relaxation play a central role in our proposed approach. The programs 9 and 12 essentially aim to learn the underlying governing function of the differential equation 5 and its relaxation. Our results show that we can satisfy these constraints over a compact set, albeit up to a margin. The existence of this margin can have non-trivial impacts. For example, small approximation errors within the value function can have larger impacts on the task cost. We hypothesise that the task cost slower convergence is a result of this phenomenon. Additionally, this margin can pose challenges in safety-critical tasks where Lyapunov guarantees are required and exact state information is not available. The

credit assignment problem, in our case, remains an issue. In our results, the HJB constraint is defined through the assigned cost function. This cost function essentially encodes the optimality constraint and, if defined poorly, can result in poor task performance. The choice of time horizon is also non trivial. In our experiments we initially perform line search over a select number of horizons and select the best one. However, our method is robust to a variety of horizons so long as values or not chosen to be very high or low. Additionally, the discretisation time Δt is also chosen by considering both simulation efficiency and approximation errors.

Finally, our method is online which can lead to typical problems associated with online algorithms such as non-stationary learning which can result in forgetting, especially in long horizon tasks.

We aim to extend this work in other further practical directions. For instance, the problem formulation and the environments used were all deterministic. This work can be extended to the stochastic domain where the benefits of sampling for example smoothing or exploration can be evaluated. Additionally, in this work, we primarily focused on problem formulation and empirical verification of the proposal. However, in further work, we aim to evaluate this method in more complex discontinuous environments with contact dynamics. However it is important to mention the ability to do this requires the availability of differentiable GPU based simulators which are currently in early stages of development. Finally, in our method, we derived the HJB policy under quadratic regularisation. This can be a limiting factor in cases where different input constraints are required. In future work, we aim to show that this extension can be relatively trivial as analytical policies are computable for any convex positive semi-definite function in u [17].

VII. CONCLUSION

We started by asking two questions. Can widely used RL algorithms such as SAC and PPO solve simple continuous control problems with minimal reward and dynamics shaping? Our results show that RL still requires significant effort in tuning to be able to achieve reasonable performance. Our second question asks for an alternative. We answer this by introducing two mathematical programs that use the HJB equation and first order gradients to learn value and Lyapunov functions. We demonstrate our effectiveness empirically by comparing to PPO and SAC on linear and nonlinear controlaffine tasks. Our results show that we can outperform both in terms of quality of the generated solution, task cost, variance in results and training stability.

REFERENCES

- [1] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter, "Learning agile and dynamic motor skills for legged robots," <u>Science Robotics</u>, vol. 4, no. 26, p. eaau5872, 2019.
- [2] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert,

- G. Powell, A. Ray et al., "Learning dexterous in-hand manipulation," The International Journal of Robotics Research, vol. 39, no. 1, pp. 3–20, 2020.
- [3] B. Katz, J. D. Carlo, and S. Kim, "Mini cheetah: A platform for pushing the limits of dynamic quadruped control," in 2019 International Conference on Robotics and Automation (ICRA), 2019, pp. 6295–6301.
- [4] Y. Tassa, T. Erez, and E. Todorov, "Synthesis and stabilization of complex behaviors through online trajectory optimization," in <u>2012 IEEE/RSJ International</u> <u>Conference on Intelligent Robots and Systems</u>. IEEE, 2012, pp. 4906–4913.
- [5] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, "Deep reinforcement learning that matters," in <u>Proceedings of the AAAI conference on artificial intelligence</u>, vol. 32, no. 1, 2018.
- [6] H. J. Suh, M. Simchowitz, K. Zhang, and R. Tedrake, "Do differentiable simulators give better policy gradients?" in <u>International Conference on Machine Learning</u>. PMLR, 2022, pp. 20 668–20 696.
- [7] S. Kuindersma, R. Deits, M. Fallon, A. Valenzuela, H. Dai, F. Permenter, T. Koolen, P. Marion, and R. Tedrake, "Optimization-based locomotion planning, estimation, and control design for the atlas humanoid robot," <u>Autonomous</u> robots, vol. 40, pp. 429–455, 2016.
- [8] I. O. Sandoval, P. Petsagkourakis, and E. A. del Rio-Chanona, "Neural odes as feedback policies for nonlinear optimal control," arXiv preprint arXiv:2210.11245, 2022.
- [9] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in <u>International</u> <u>conference on machine learning</u>. PMLR, 2018, pp. 1861– 1870.
- [10] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," arXiv preprint arXiv:1707.06347, 2017.
- [11] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in <u>International conference on machine learning</u>. PMLR, 2015, pp. 1889–1897.
- [12] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," <u>arXiv:1509.02971</u>, 2015.

- [13] B. Recht, "A tour of reinforcement learning: The view from continuous control," Annual Review of Control, Robotics, and Autonomous Systems, vol. 2, pp. 253–279, 2019
- [14] I. Mordatch, K. Lowrey, G. Andrew, Z. Popovic, and E. V. Todorov, "Interactive control of diverse complex characters with neural networks," <u>Advances in neural</u> information processing systems, vol. 28, 2015.
- [15] J. Wang, T. S. Lembono, S. Kim, S. Calinon, S. Vi-jayakumar, and S. Tonneau, "Learning to guide online multi-contact receding horizon planning," in 2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2022, pp. 12 942–12 949.
- [16] H. Li, R. J. Frei, and P. M. Wensing, "Model hierarchy predictive control of robotic systems," <u>IEEE Robotics and</u> Automation Letters, vol. 6, no. 2, pp. 3373–3380, 2021.
- [17] M. Lutter, S. Mannor, J. Peters, D. Fox, and A. Garg, "Value iteration in continuous actions, states and time," arXiv preprint arXiv:2105.04682, 2021.
- [18] K. Lowrey, A. Rajeswaran, S. Kakade, E. Todorov, and I. Mordatch, "Plan online, learn offline: Efficient learning and exploration via model-based control," <u>arXiv preprint</u> arXiv:1811.01848, 2018.
- [19] C. Dawson, Z. Qin, S. Gao, and C. Fan, "Safe nonlinear control using robust neural lyapunov-barrier functions," in <u>Conference on Robot Learning</u>. PMLR, 2022, pp. 1724–1735.
- [20] W. Xiao, T.-H. Wang, R. Hasani, M. Chahine, A. Amini, X. Li, and D. Rus, "Barriernet: Differentiable control barrier functions for learning of safe robot control," <u>IEEE</u> Transactions on Robotics, pp. 1–19, 2023.
- [21] Y.-C. Chang, N. Roohi, and S. Gao, "Neural lyapunov control," <u>Advances in neural information processing systems</u>, vol. 32, 2019.
- [22] W. Xiao, R. Hasani, X. Li, and D. Rus, "BarrierNet: A Safety-Guaranteed Layer for Neural Networks," <u>arXiv</u> e-prints, p. arXiv:2111.11277, Nov. 2021.
- [23] S. Ainsworth, K. Lowrey, J. Thickstun, Z. Harchaoui, and S. Srinivasa, "Faster policy learning with continuoustime gradients," in <u>Learning for Dynamics and Control</u>. PMLR, 2021, pp. 1054–1067.
- [24] X. Zhang, J. Long, W. Hu, W. E, and J. Han, "Initial value problem enhanced sampling for closed-loop optimal control design with deep neural networks,"

- 2023. [Online]. Available: https://openreview.net/forum? id=oXM5kdnAUNq
- [25] S. Bansal and C. J. Tomlin, "Deepreach: A deep learning approach to high-dimensional reachability," in 2021 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2021, pp. 1817–1824.
- [26] S. Engin and V. Isler, "Neural optimal control using learned system dynamics," <u>arXiv preprint</u> arXiv:2302.09846, 2023.
- [27] Y. D. Zhong, B. Dey, and A. Chakraborty, "Symplectic ode-net: Learning hamiltonian dynamics with control," in <u>International Conference on Learning Representations</u>, 2020. [Online]. Available: https://openreview.net/forum? id=ryxmb1rKDS
- [28] D. H. Jacobson, "New second-order and first-order algorithms for determining optimal control: A differential dynamic programming approach," <u>Journal of Optimization</u> Theory and Applications, vol. 2, pp. 411–440, 1968.
- [29] D. Liberzon, <u>Calculus of variations and optimal control</u> theory: a concise introduction. Princeton university press, 2011.
- [30] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud, "Neural ordinary differential equations," <u>Advances</u> in neural information processing systems, vol. 31, 2018.
- [31] Y. Lin and E. D. Sontag, "A universal formula for stabilization with bounded controls," <u>Systems & control</u> letters, vol. 16, no. 6, pp. 393–397, 1991.
- [32] B. Amos, L. Xu, and J. Z. Kolter, "Input convex neural networks," in Proceedings of the 34th International Conference on Machine Learning, ser. Proceedings of Machine Learning Research, vol. 70. PMLR, 2017, pp. 146–155.
- [33] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann, "Stable-baselines3: Reliable reinforcement learning implementations," <u>Journal of Machine Learning Research</u>, vol. 22, no. 268, pp. 1–8, 2021. [Online]. Available: http://jmlr.org/papers/v22/20-1364.html
- [34] R. Tedrake, <u>Underactuated Robotics</u>, 2023. [Online]. Available: https://underactuated.csail.mit.edu
- [35] E. Todorov and W. Li, "Optimal control methods suitable for biomechanical systems," in <u>Proceedings of the 25th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (IEEE Cat.</u>

- No.03CH37439), vol. 2, 2003, pp. 1758-1761 Vol.2.
- [36] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," 2016.
- [37] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," arXiv preprint arXiv:1412.6980, 2014.