
DARE: The Deep Adaptive Regulator for Control of Uncertain Continuous-Time Systems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A fundamental challenge in continuous-time optimal control (OC) is the efficient
2 computation of adaptive policies when agents act in unknown, uncertain environ-
3 ments. Traditional OC methods, such as dynamic programming, face challenges in
4 scalability and adaptability due to the curse-of-dimensionality and the reliance on
5 fixed models of the environment. One approach to address these issues is Model
6 Predictive Control (MPC), which iteratively computes open-loop controls over
7 a receding horizon. However, classical MPC algorithms typically also assume
8 a fixed environment. Another approach is Reinforcement Learning (RL) which
9 scales well to high-dimensional setups but is often sample inefficient. Certain RL
10 methods can also be unreliable in highly stochastic continuous-time setups and
11 may be unable to generalize to unseen environments. This paper presents the **Deep**
12 **Adaptive Regulator (DARE)** which uses physics-informed neural network based
13 approximations to the agent’s value function and policy which are trained online
14 to adapt to unknown environments. To manage uncertainty of the environment,
15 DARE optimizes an augmented reward objective which dynamically trades off ex-
16 ploration with exploitation. We show that our method effectively adapts to unseen
17 environments in settings where “classical” RL fails and is suited for online adaptive
18 decision-making in environments that change in real time.

19 1 Introduction

20 Many real-world decision making problems in fields such as biology [24] and algorithmic trading [11],
21 require agents to act at high-frequency in noisy systems, and hence can be modeled as continuous-time
22 stochastic optimal control (OC) problems. Moreover, in many of these areas, the agent is uncertain of
23 the system’s true dynamics, which may be non-stationary, and the agent must learn these dynamics
24 in real-time through interacting with the environment [17]. However, classical continuous-time OC
25 methods for generating optimal policies, such as solving the standard Hamilton-Jacobi-Bellman
26 (HJB) equation, generally assume known, fixed environments, limiting their wider application [55].
27 To address this challenge, *adaptive* control methods enable agents to optimize their performance by
28 continuously learning and adapting to the unknown and evolving dynamics of their environment [8].
29 However, these methods typically rely on explicit parametric models of the environment [12], and the
30 numerical methods required to solve these problems are often intractable [10].

31 One common approach to controlling *unknown* systems is Reinforcement Learning (RL), in which
32 agents compute (near) optimal policies iteratively through trial and error [48]. However, without
33 certain ad-hoc techniques such as action-repetition, certain RL methods can perform poorly in
34 highly stochastic and nearly continuous-time setups, even when the environment is stationary; see
35 [49, 54, 25]. Moreover, model-free RL methods suffer from sample inefficiency [3, 13] which can
36 limit their applicability to real-world problems, whereas model-based methods that aim to mitigate

37 this drawback can require expert data that is not always accessible and are likely to lead to poor
 38 generalization [39].

39 Another control methodology for the control of uncertain systems is model predictive control (MPC).
 40 In MPC, agents optimize control inputs by solving a sequence of *open-loop* optimization problems
 41 based on a predictive model over a receding time horizon, allowing for real-time adjustments of
 42 the agent’s policy; see [19]. While MPC is generally more computationally efficient than standard
 43 dynamic programming methods, which often require compute-intensive finite-difference schemes,
 44 standard methods of solving MPC problems still face computational bottlenecks, hindering their
 45 application to high-frequency environments [51]. For a more detailed review of relevant literature,
 46 see Appendix A. Hence, flexible methods that compute optimal policies efficiently in uncertain,
 47 non-stationary environments are desirable.

48 In this paper, we propose the Deep Adaptive Regulator (DARE). DARE is a deep learning-based
 49 method for solving stochastic continuous-time adaptive control problems. Our method consists of
 50 two distinct phases: offline and online. In both phases, physics-informed neural networks (PINNs)
 51 [42] parameterize the agent’s value function and policy. During both phases, we optimize the value
 52 function and policy networks to satisfy the HJB equation resulting from the agent’s estimate of the
 53 dynamics of the system and costs.

54 In the *offline* phase, the agent is endowed with an initial (potentially misspecified) estimate of the
 55 environment, which in practice may be learned from historical data. With this initial estimate, the
 56 agent constructs and solve an approximate OC problem which yields an initial policy. To accelerate
 57 training, we make use of two inductive biases: 1) we let the value function network learn the
 58 discrepancy between the true value function and the terminal condition of the HJB equation, and 2)
 59 we let the policy network learn the discrepancy between the true policy and a locally optimal policy
 60 generated by the Iterative Linear-Quadratic Gaussian (ILQG) method [51]. In the *online* phase, the
 61 agent implements their policy in the true environment and updates their estimate of the environment
 62 based on (potentially noisy) observations of the system and running costs. At each observation time,
 63 the agent updates their value function and policy by solving an updated OC problem over a receding
 64 horizon, similar to MPC. In contrast to MPC, which iteratively computes a sequence of open-loop
 65 controls to attain closed-loop feedback, DARE approximates the optimal closed-loop policy directly
 66 and updates this approximation in real-time. Moreover, to account for the agent’s uncertainty about
 67 the environment, we consider a modification of the agent’s objective function which dynamically
 68 balances exploration and exploitation.

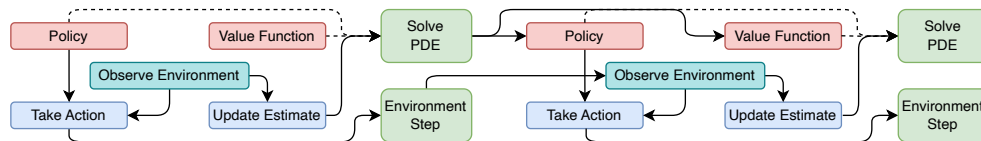


Figure 1: A schematic of DARE.

69 To benchmark DARE, we study its performance in three adaptive OC problems: (i) a one-dimensional
 70 Linear-Quadratic-Gaussian (LQG) Regulator problem in which the agent is uncertain of the drift of
 71 the system, (ii) an adaptive nonlinear MPC problem in which the agent is uncertain of the running
 72 cost and models it in real time with a Gaussian Process (GP), and (iii) a realistic high-dimensional
 73 nonlinear MPC problem motivated by algorithmic trading in finance [11]. While these problems are
 74 relatively simple, we demonstrate that in the LQG problem, a suite of classical RL methods including
 75 PPO [44], A2C [35], and SAC [20] are unable to learn robust policies when the time between actions
 76 is small, for a variety of action-repetition parameters. Moreover, we demonstrate that when RL
 77 methods are able to learn a robust solution, they require substantially more training time than DARE.
 78 In contrast, we show that DARE can learn accurate solutions and adapt to non-stationary environments
 79 in each problem.

80 In summary, this paper: (i) proposes the PINN-based method DARE to compute adaptive control
 81 policies in uncertain, continuous-time environments (ii) proposes an OC problem formulation that
 82 explicitly trades off exploration and exploitation which aids adaptation in non-stationary systems,
 83 (iii) demonstrates experimentally that classical RL methods fail on relatively simple continuous-time
 84 adaptive control tasks when time between actions is small, and (iv) proposes a inductive bias based

85 on ILQG and the structure of the HJB which significantly accelerates offline training and improves
 86 speed of adaptation in non-stationary environments, outperforming classical RL.

87 2 Problem Formulation

88 Let $X_t \in \mathbb{R}^{d_x}$ be a stochastic system evolving continuously in time. We consider an agent who
 89 controls X with a policy $u_t \in \mathbb{R}^{d_u}$ over a fixed time horizon $T > 0$ to maximize a terminal reward
 90 $g : \mathbb{R}^{d_x} \rightarrow \mathbb{R}$. The agent's actions u_t on the system X_t incur a penalty modelled by a function
 91 $f : \mathbb{R}^{d_x} \times \mathbb{R}^{d_u} \rightarrow \mathbb{R}$, and their impact on the system dynamics is modelled by a drift function
 92 $h : \mathbb{R}^{d_x} \times \mathbb{R}^{d_u}$. The system evolves according to the dynamics

$$dX_t = h(X_t, u_t) dt + \tilde{\Sigma} dW_t, X_0 \in \mathbb{R}^{d_x}, \quad (1)$$

93 where W is a d_x -dimensional Brownian motion and $\tilde{\Sigma} \in \mathbb{R}^{d_x \times d_x}$ is a covariance matrix. We assume
 94 $\tilde{\Sigma}, T$ and g are fixed and known to the agent, and $\mathfrak{p} := (h, f)$ represents the modelling assumptions
 95 of the agent over the environment. We refer to \mathfrak{p} as the *OC pair*.

96 Classical OC approaches assume a fixed and known pair \mathfrak{p} to compute an optimal policy. In practice,
 97 the agent uses an uncertain estimate $\hat{\mathfrak{p}}$ of the true environment. To account for this uncertainty, DARE
 98 solves the decision-making problem in two phases: offline and online. In the offline phase, the agent
 99 solves an OC problem according to an initial estimate of the environment $\hat{\mathfrak{p}}_0$. In the online phase, the
 100 agent implements the policy from the offline phase in the true environment, receives noisy samples
 101 of the true OC pair, and updates their estimate of the environment. Then, the agent updates their
 102 control policy to be the solution of an updated OC problem according to the agent's new estimate of
 103 the environment. The agent is uncertain of their estimate, so it may be profitable to explore unknown
 104 areas of the system for potentially higher rewards. Hence, a balance must be struck between exploring
 105 new information and exploiting existing knowledge.

106 **Offline phase.** At time $t = 0$, the agent assumes an initial estimate $\hat{\mathfrak{p}}_0 = (\hat{h}_0, \hat{f}_0)$ of the OC pair.
 107 To explicitly account for the exploration-exploitation trade-off, the agent seeks an optimal policy u^*
 108 which maximizes the following performance criterion:

$$J(s, x; u) = \mathbb{E} \left[g(X_T) - \int_s^T \mathbb{E} [\hat{f}_0] (X_r, u_r) dr + \phi \int_s^T \text{Var} [\hat{\mathfrak{p}}_0] (X_r, u(r, X_r)) dr \mid \mathcal{G}_s \right], \quad (2)$$

109 for all $s \in [0, T]$, where \mathcal{G}_s is the information known to the agent at time s , $\mathbb{E} [\hat{f}_0]$ denotes the mean
 110 prediction of the estimate \hat{f}_0 and $\text{Var} [\hat{\mathfrak{p}}_0]$ is the sum of the variance of each estimator in $\hat{\mathfrak{p}}_0$,
 111 and we assume that X_s follows the dynamics

$$dX_s = \hat{h}_0(X_s, u_s) ds + \tilde{\Sigma} dW_s, \quad X_0 \in \mathbb{R}^{d_x}. \quad (3)$$

112 The objective optimized by the agent in (2) is an adjusted formulation of the agent's true objective,
 113 in which the agent explicitly rewards or penalizes uncertainty on their estimate of the environment.
 114 More precisely, When $\phi > 0$ (resp. < 0) the agent rewards (resp. penalizes) exploration, i.e., the
 115 agent is encouraged to visit areas of the environment with higher (resp. lower) uncertainty. We
 116 show in Section 4.5 that the exploration parameter ϕ is key to the performance of decision-making
 117 problems in noisy and non-stationary environments.

118 The incorporation of the variance of the environment estimation in (2) is similar in spirit to the
 119 variance adjusted objective common in bandit algorithms such as [46]. However a notable difference
 120 is that the uncertainty in (2) is incorporated *dynamically* as opposed to myopically in the case of
 121 bandits.

122 To solve the problem (2), the agent defines the value function

$$V(s, x) = \sup_u J(s, x; u). \quad (4)$$

123 We assume that the dynamic programming principle holds for $\mathbb{E} [\hat{f}_0]$ and $\text{Var} [\hat{\mathfrak{p}}_0]$, so V solves the
 124 HJB equation:

$$0 = V_t + \frac{1}{2} \text{Tr}(\Sigma \nabla_{xx} V) + \sup_{u \in \mathbb{R}^{d_u}} H(x, u, \nabla_x V(t, x); \hat{\mathfrak{p}}_0), \quad (5)$$

125 subject to terminal condition $V(T, x) = g(x)$, where $\Sigma = \tilde{\Sigma} \tilde{\Sigma}^\top$. For $\ell \in \mathbb{R}^{d_x}$, the Hamiltonian H
 126 in (5) is defined as

$$H(x, u, \ell; \hat{\mathbf{p}}_0) = \hat{h}_0(x, u)^\top \ell + \mathbb{E} \left[\hat{f}_0 \right] (x, u) - \phi \text{Var} \left[\hat{\mathbf{p}}_0 \right] (x, u(t, x)). \quad (6)$$

127 The policy u^* which maximizes (2) can be computed using the following first order conditions for
 128 $s \in [0, T]$:

$$u^*(s, x; \hat{\mathbf{p}}_0) = \arg \max_{u \in \mathbb{R}} H(x, u, V_x(s, x); \hat{\mathbf{p}}_0), \quad (7)$$

129 where $V(t, x)$ solves the nonlinear PDE (5). In contrast to several approaches in RL which address
 130 the exploration-exploitation trade-off through penalization or reward of *random* control processes (see
 131 [53] in a continuous-time setup), our method learns a control policy that is a *deterministic* function of
 132 the environment and which explores domain regions in which the agent is uncertain of their estimates
 133 of the OC pair.

134 **Online Problem.** At each time $t \in (0, T]$, the agent takes an action u_t , observes a noisy sample
 135 of the true environment $\mathbf{p}(u_t, X_t) + \epsilon_t$ for some i.i.d. noise $\{\epsilon_t\}$, and updates their estimate $\hat{\mathbf{p}}_t$
 136 accordingly.¹ Conditionally on the new estimate, the policy computed during the offline phase is not
 137 optimal. To adapt the optimal policy, the agent computes a new policy which maximizes the updated
 138 objective, for $s \in [t, T]$:

$$J(s, x; u) = \mathbb{E} \left[g(X_T) - \int_s^T \mathbb{E} \left[\hat{f}_t \right] (X_r, u_r) dr + \phi \int_s^T \text{Var} \left[\hat{\mathbf{p}}_t \right] (X_r, u(r, X_r)) dr \mid \mathcal{G}_s \right]. \quad (8)$$

139 This new control policy is then implented in the true system for $s \in [t, T]$ un-
 140 til the agent observes the system again, at which point the procedure is repeated.
 141 When the uncertainty $\text{Var}(\hat{\mathbf{p}}_t)$ is high and the
 142 agent rewards exploration, i.e., $\phi < 0$, the DARE
 143 policy focuses on improving the agent’s estimate
 144 of their environment. As the estimation
 145 accuracy increases, the variance of the estima-
 146 tor and hence its contribution to the objective
 147 decreases. Hence, the DARE policy naturally
 148 balances the exploration-exploitation trade-off
 149 throughout the online phase. We demonstrate
 150 the benefit of the exploration term to overcome
 151 misspecified priors or nonstationary environ-
 152 ments in Section 4.1. We note that simply opti-
 153 mizing the offline variance-adjusted objective
 154 (2) without updating the system and re-solving
 155 for a new optimal control would likely lead to
 156 sub-optimal solutions, as the variance penalty
 157 would be fixed for the entire time horizon. As
 158 far as we are aware, incorporation of model un-
 159 certainty into continuous-time stochastic control
 160 problems in this fashion is new to the literature.

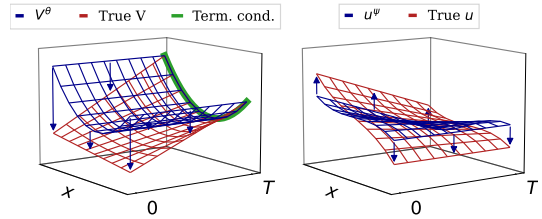


Figure 2: Illustration of the initialization of V^{θ_0}, u^{ψ_0} in the offline phase. The value function network is initialized around the terminal condition g in (2), and the control policy network is initialized around an affine approximation to the true optimal control.

161 3 The Deep Adaptive Regulator

162 A clear difficulty faced in the procedure above is the computation of the updated policy in the online
 163 phase. Even in low-dimensional settings, compute-intensive finite-difference schemes are required
 164 to solve the resulting HJBs [18]. Hence, we present DARE, which addresses this issue with PINN
 165 models of both the value function and control policy.

166 **Offline Phase.** To obtain the initial policy corresponding to the pair $\hat{\mathbf{p}}_0$, we use PINN approxi-
 167 mations V^θ and u^ψ of the value function V and the optimal control u that solve the HJB (5). We

¹We do not assume a particular estimation procedure, but this can be achieved with function approximators suitable for online learning, e.g., Gaussian Processes or Bayesian NNs as in [16].

168 initialize V^θ and u^ψ , for $s \in [0, T]$, as follows:

$$\begin{cases} V^\theta(s, x) &= g(x) + \mathfrak{X}^\theta(s, x), \\ u^\psi(s, x) &= \hat{u}_{X_0}(s, x) + \mathfrak{X}^\psi(s, x), \end{cases} \quad (9)$$

169 where g is the terminal reward function and \hat{u}_{X_0} is a locally optimal linear approximation of the true
170 optimal control. We use the ILQG method of [51] to compute \hat{u}_{X_0} starting from X_0 and we set \mathfrak{X}^θ
171 and \mathfrak{X}^ψ to be fully-connected feedforward networks; see Figure 2.

172 To train V^θ and u^ψ , we devise a multi-objective PINN loss which considers (i) the HJB (5), (ii) the
173 Hamiltonian (6) satisfying first-order conditions, and (iii) the terminal condition. We use Monte
174 Carlo integration to minimize the loss function on a compact domain $K \subset \mathbb{R}^{dx}$. We let $\|\cdot\| =$
175 $\|\cdot\|_{L^2([0, T] \times K)}$ and we reformulate (5) in a variational form to define the loss:

$$\mathcal{L}(\theta, \psi; \hat{\mathbf{p}}) = \mathcal{L}_{\text{HJB}} + \mathcal{L}_{\text{hamiltonian}} + \mathcal{L}_{\text{terminal}}, \quad (10)$$

176 where

$$\begin{cases} \mathcal{L}_{\text{HJB}} = \|V_t^\theta + \frac{1}{2} \text{Tr}(\Sigma \nabla_{xx} V^\theta) + H(\cdot, u^\psi(\cdot, \cdot), V_x^\theta(\cdot, \cdot); \hat{\mathbf{p}})\|, \\ \mathcal{L}_{\text{hamiltonian}} = \|\partial_u H(\cdot, u^\psi(\cdot, \cdot), V_x^\theta(\cdot, \cdot); \hat{\mathbf{p}})\|, \\ \mathcal{L}_{\text{terminal}} = \|V^\theta(T, \cdot) - g\|. \end{cases} \quad (11)$$

177 The loss (10) is similar to that in [1], however, here we use a first-order condition for the Hamiltonian
178 component, which we found to improve training performance when H is concave. Otherwise, we set

$$\mathcal{L}_{\text{hamiltonian}} = -\|H(\cdot, \cdot, u^\psi; \hat{\mathbf{p}})\|. \quad (12)$$

179 We summarize the procedure in Algorithm B.

180 **Online Phase.** Let $\mathcal{T} = \{t_0, \dots, t_n\} \subset [0, T]$ be a set of potentially irregularly spaced times which
181 are unknown to the agent at time $t = 0$. In the online phase, the agent uses new estimates of the
182 environment to update their control policy as follows. Suppose the agent calculated $V^{\theta_{t_{k-1}}}(\cdot, \cdot; \hat{\mathbf{p}}_{t_{k-1}})$
183 and $u^{\psi_{t_{k-1}}}(\cdot, \cdot; \hat{\mathbf{p}}_{t_{k-1}})$ at time t_k , where $\theta_{t_{k-1}}$ and $\psi_{t_{k-1}}$ minimize the loss $\mathcal{L}(\theta, \psi, \hat{\mathbf{p}}_{t_{k-1}})$. At time
184 t_k , the agent (i) takes the action $u^{\psi_{t_{k-1}}}(t_k, X_{t_k}; \hat{\mathbf{p}}_{t_{k-1}})$ and (ii) computes the new estimate $\hat{\mathbf{p}}_{t_k}$. Then,
185 over the period $[t_k, t_{k+1})$, the agent minimizes $\mathcal{L}(\theta, \psi, \hat{\mathbf{p}}_{t_k})$ to compute the parameters $(\theta_{t_k}, \psi_{t_k})$.
186 This loss minimization uses a gradient-based method (e.g., ADAM), with $(\theta_{t_{k-1}}, \psi_{t_{k-1}})$ as a warm
187 start for the neural networks; our method is outlined in Figure 1 and Algorithm B.

188 4 Numerical Experiments

189 This section investigates the performance of DARE in both the offline and online phases. Three
190 control tasks are considered, including Linear-Quadratic Gaussian (LQG) control, an augmentation
191 of the LQG problem with uncertain running costs, and a high-dimensional control task motivated by
192 algorithmic trading in finance. First, we demonstrate the ability of DARE to learn accurate solutions
193 in the offline phase. We then present a failure mode of classical RL methods in continuous-time
194 control tasks, demonstrating degradation of performance in the simple LQG control problem. We
195 also investigate the impact of the variance-adjusted objective in the MPC problem when running cost
196 observations are corrupted by noise. Finally, we present a comparison of the sample efficiency of
197 DARE compared to RL methods, demonstrating vastly improved training times.

198 4.1 Description of Control Tasks

199 **Linear-Quadratic-Gaussian.** Consider system dynamics

$$dX_t = (b + c u_t) dt + \tilde{\Sigma} dW_t, \quad X_0 \in \mathbb{R}, \quad (13)$$

200 where b is a constant drift, $c > 0$ scales the linear impact of an agent on the system, and $\tilde{\Sigma} > 0$ is the
201 variance of the observation noise. The agent maximizes the LQ criterion

$$\mathbb{E} \left[X_T - \alpha X_T^2 - \phi \int_0^T u_t^2 dt \right], \quad (14)$$

202 where $\phi > 0$ scales the running quadratic penalty and $\alpha > 0$ scales the terminal quadratic penalty.

203 **Model Predictive Control.** We consider an augmentation of LQG in which the agent is subject
 204 unknown running cost and must model this cost in real-time. In particular, the dynamics of the system
 205 evolves according to (13), and the true running cost of the system is quadratic as in (14). The agent
 206 uses a Gaussian Process (GP) \hat{f} to model the true running cost, and updates this model in real time
 207 using noisy observations of accumulated running costs. Hence, in both the offline and online phase
 208 of training, the DARE optimizes the objective

$$\mathbb{E} \left[X_T - \alpha X_T^2 - \phi \int_0^T \mathbb{E}[\hat{f}](u_t) dt - \varphi \int_0^T \text{Var}[\hat{f}](u_t) dt \right]. \quad (15)$$

209 **High-Dimensional Control.** For a high-dimensional control example, we consider an example
 210 from finance. In recent years, regulators have urged financial institutions to manage the risk of their
 211 trading activity within very large portfolios called central risk books. The aggregated trading activity
 212 of large institutions is often conducted at very high frequency and can be modeled as an OC problem.
 213 The controlled system is described by the agent’s inventory $Q_t \in \mathbb{R}^d$, the asset prices $S_t \in \mathbb{R}^d$, and
 214 running wealth $X_t \in \mathbb{R}$, with dynamics:

$$dQ_t = u_t dt, \quad dS_t = \tilde{\Sigma} dW_t, \quad dX_t = -u_t^\top S_t dt - f(u_t) dt,$$

215 where $u_t \in \mathbb{R}^d$ denotes the trader’s speed of trading. The agent incurs transaction costs according to
 216 some unknown function of the trading speed $f(u_t) \in \mathbb{R}^d$, and maximizes the exponential utility of
 217 their terminal wealth for some estimate \hat{f} of the true transaction costs

$$\sup_v \mathbb{E} \left[-\exp(-\gamma (X_T + Q_T^\top S_T - Q_T^\top \Gamma Q_T)) \right],$$

218 In this example, we set $d = 5$. In Appendix H we record the model parameters η and $\tilde{\Sigma}$ and include
 219 a detailed motivation for this problem.

220 4.2 Offline Performance

221 First, we investigate the performance of DARE in the offline phase. To do so, we fix the agent’s prior,
 222 and demonstrate the ability of DARE to learn accurate solutions to the HJB equation posed with this
 223 prior. In particular, we investigate the impact of the inductive biases used to augment the value
 224 function and control policy defined in (9). We compare DARE to PINNs without such bias, that is, we
 225 consider two methods, MLP and DGM with initializations

$$V^\theta(t, x) = \mathfrak{X}^\theta(t, x) \quad \text{and} \quad u^\psi(t, x) = \mathfrak{X}^\psi(t, x),$$

226 where \mathfrak{X}^θ and \mathfrak{X}^ψ are feedforward DNNs with Xavier initialization in MLP and LSTM-like networks
 227 as in the Deep Galerkin Method [45] in DGM.

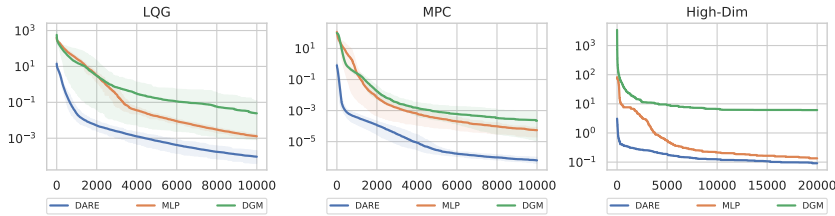


Figure 3: We plot the training loss trajectory in the offline phase for DARE, MLP, and DGM. In the MPC problem, we set $b = 0$, and other model parameters are given in Appendix C. In the LQG and MPC figures, we take the average of 100 seeds.

228 In each of problems, the MLPs of both DARE and MLP, there are 2 layers and 20 hidden units. In DGM,
 229 there are two hidden LSTM-like layers between two single layer feedforward neural networks of
 230 width 20. Additional details about the parameters for each problem in this experiment are found in
 231 Appendix C. Figure 3 shows the loss (10) of the three methods throughout offline training when the
 232 agent’s prior is fixed. On average, we observe that DARE substantially outperforms other architectures
 233 in convergence speed in all three problems, achieving lower loss with fewer iterations. While existing
 234 work [45] emphasizes the importance of network architecture for performance, our findings indicate
 235 that, at least for simple problems, inductive biases may hold greater importance.

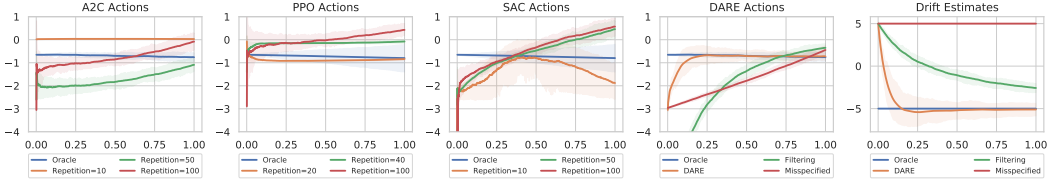


Figure 5: We plot the actions each of algorithm in the online phase, along with the drift estimation used by the RL agents and DARE. For each RL method, we plot actions corresponding to the top three performing methods in Figure 4. We also include continuous-time Kalman filtering using a prior $b \sim \mathcal{N}(5, 3)$, which is calculated using techniques in Appendix G.

236 4.3 Failure Mode of Reinforcement Learning in Continuous-Time Systems

237 In [49], the authors find that Q -learning
 238 based RL methods struggle in near
 239 continuous-time environments, as the Q -
 240 value contribution of a single action
 241 vanishes with a shrinking discretization
 242 timestep. Here, we demonstrate that a suite
 243 of RL methods, including Proximal Policy
 244 Optimization (PPO) [44], Advantage
 245 Actor Critic (A2C) [35], and Soft Actor
 246 Critic (SAC) [20] struggle in even the most
 247 simple continuous-time stochastic control
 248 setting, the LQG problem, when the time
 249 discretization is very small.

250 In the LQG problem the drift b is unknown
 251 to the agent. To train the RL agent, we
 252 simulate trajectories of the state process
 253 in (13) using samples of the drift b from
 254 some prior distribution. We then let the
 255 drift estimation be a state variable, that is,
 256 we let the policy of each agent $u^\psi = u^\psi(t, x; b)$. To estimate the drift, at each observation time t
 257 in the online phase, each agent uses an exponential moving average with smoothing $\lambda = 0.95$
 258 to determine an estimate b_t . To test each agent’s ability to adapt, we consider the case of a misspecified
 259 prior distribution where the drift is sampled from $b \sim \mathcal{N}(b_0, \Pi_0)$, where $b_0 = 5$ and $\Pi_0 = 3$. To
 260 implement DARE in this environment, we use the same drift estimation, and after each observation,
 261 we use 10 ADAM steps to update the DARE policy. Additional details regarding our implementation
 262 and hyperparameters used for A2C, PPO, and SAC are included in Appendix J.

263 In Figure 4, we plot the performance of each RL agent when the system is simulated with the true
 264 drift $b = 5$ and the agent must estimate the drift in run time. We discretize the problem horizon into
 265 1000 steps and investigate the effect of varying number of times actions are repeated. When actions
 266 are repeated less frequently, we observe that the performance of each RL algorithm steadily degrades.
 267 DARE, however, does not require action-repetition during training, and achieves similar performance
 268 to the best performing RL method, which required extensive hyperparameter tuning.

269 In Figure 5, we picked the top three performing action-repetition hyperparameters for each RL
 270 method and plotted the actions of these agents in the online phase. Note that while A2C with 100
 271 repeated actions is the top performing method, the actions produced by this method are quite far
 272 from the oracle. On the other hand, DARE learns the oracle control policy quickly after only a few
 273 environment steps.

274 4.4 Real Time Adaptation in Nonlinear Environments

275 Here, we investigate the performance of the online phase of DARE in the MPC and high-dimensional
 276 control problems. To test the ability of DARE to adapt in these settings, we first train DARE to a

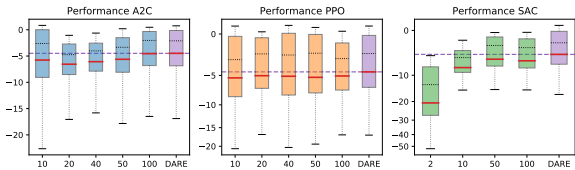


Figure 4: Average total reward of PPO, A2C, SAC, and DARE in the LQG problem across 1000 rollouts. The oracle average reward is denoted by the dashed line and is calculated using techniques outlined in Appendix G. Each agent is trained offline using a misspecified prior of $b_0 = 5$ in the case of DARE and $b \sim \mathcal{N}(b_0, \Pi_0)$ for the RL agents, and upon rollout, uses an exponential moving average to estimate the true drift $b \sim \mathcal{N}(-5, 3)$ of the system.

277 misspecified prior in the offline phase. Then, we suppose the agent learns the true environment
 278 parameters, and we test how many iterations are required to adjust the misspecified policy to the
 279 optimal policy of the true environment.

280 For the MPC problem, We define two running penalty functions $f_i = |u|^{1+\gamma_i}$ for $i \in \{0, 1\}$, where
 281 $\gamma_0 = 1.3$ and $\gamma_1 = 1$. We consider agents who use a Gaussian Process (GP) \hat{f} as a predictive model
 282 for the running penalty; see Appendix I for details on GPs. First, the agents fit two Gaussian Process
 283 \hat{f}_i for $i \in \{0, 1\}$ to ten noisy, random samples of the running penalty f_i . Next, we use DARE, MLP,
 284 and DGM to solve for solutions V^{θ_0}, u^{ψ_0} relative to \hat{f}_0 . Once all methods have converged, we change
 285 the agents' estimate of the running penalty to \hat{f}_1 and re-train V^{θ_0}, u^{ψ_0} with this updated penalty
 286 function. For the high-dimensional control problem, we first train each algorithm using a running
 287 cost of $f(u)u^{\gamma_1} \eta u^{\gamma_0}$ where the exponent $\gamma = 1.3$ is applied element-wise, and then we change the
 288 exponent to $\gamma = 1$.

289 In Figure 6, we plot the training loss after adjusting the running costs. In the MPC problem, all
 290 methods learn the new policy with comparable precision after a few hundred iterations. In the
 291 high-dimensional problem, DGM is unable to adapt to the new environment. DARE, however, is able to
 292 adapt to the new environment faster than the other methods, requiring less than 20 ADAM steps to
 293 achieve satisfactory precision. Each iteration lasts 0.00446 seconds on average in our experiments so
 294 DARE is suited for online problems with near continuous observations in nonstationary environments.

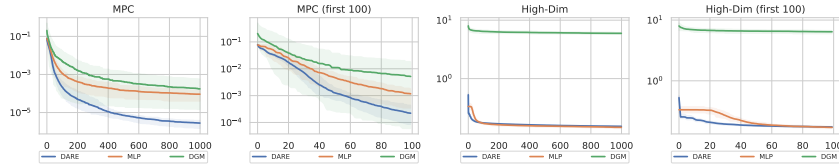


Figure 6: We plot the training loss trajectory in the online phase DARE, MLP, and DGM when the prior in each case is $\gamma = 1.3$ and the true environment is $\gamma = 1$.

295 4.5 Online Performance: Exploration-Exploitation in Non-Stationary Environments

296 We now focus on the MPC problem above. We consider the online phase and investigate when
 297 the agent's real-time observations of the running penalty are corrupted by noise. That is, the agent
 298 observes $|u^{\psi}(t, X_t)|^2 + \epsilon_t$ for $\epsilon_t \sim \mathcal{N}(0, .02)$. In this setting, we investigate whether it is beneficial
 299 for the agent to explore to ensure they have an accurate model of the running penalty function. We
 300 examine the effect of varying the variance penalty weight φ in (15). In particular, we test DARE with
 301 $\varphi = 0$ and $\varphi < 0$, which corresponds to an agent which is indifferent to exploration and encourages
 302 exploration, respectively.

303 Figure 7 shows that in the presence of noise, encouraging exploration in the objective enables the agent to learn the true policy far
 304 faster than being indifferent to exploration. Often, the absence of exploration in noisy environments leads to local optima in the value
 305 function and control policy network parameters, because a (wrong)
 306 exploration in noisy environments leads to local optima in the value
 307 function and control policy network parameters, because a (wrong)
 308 mean prediction leads to a specific policy which prevents accurate
 309 learning of the cost function in the whole domain of controls.

310 Finally, we consider a simulation setup where the true form of
 311 the cost function randomly switches between that of $\gamma_1^* = 1.3$
 312 and $\gamma_2^* = 1$ according to a Poisson process with intensity 0.005,
 313 i.e., with 1.5 switches, on average, per simulation. We consider
 314 an environment which starts with the cost functional $\gamma_1^* = 1.3$,
 315 and the observations of the running penalty $|u^{\psi}(t, X_t)|^{1+\gamma^*} + \epsilon_t$
 316 are corrupted with noise $\epsilon_t \sim \mathcal{N}(0, .1)$. Similar to the previous
 317 experiment, the agent uses a GP to model the running penalty. In
 318 Figure 8 one sees that DARE quickly adapts to the new environment
 319 after each jump in the running cost exponent.

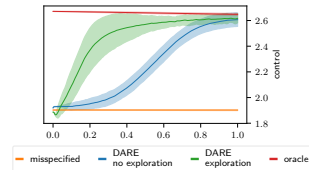


Figure 7: Mean and std dev of policy for oracle, misspecified, and DARE when $\varphi = 0$ (no exploration), and when $\varphi = 5 \cdot 10^{-3}$ (exploration).

320 We note that our decision to exclude an RL benchmark in this case
 321 was informed by two factors. First, we demonstrate above that
 322 the RL algorithms considered struggle even in the LQG setting,
 323 suggesting poor performance in nonlinear problems. Second, in the
 324 MPC problem, the agent uses a GP to model uncertainty. In the
 325 LQG problem, uncertainty about the drift was parametric and could
 326 be incorporated as a state variable for the RL agent. Hence, the
 327 RL agent could be trained on samples of different drifts and learn
 328 the optimal policy for all drifts in the sampled region. In the MPC
 329 problem, the GP estimate could not be incorporated directly as a
 330 state variable for the RL agent, because the GP is a non-parametric
 331 function approximator. Tailoring an RL algorithm to optimize over
 332 a set of non-parametric cost functions was out of the scope of this
 333 paper.

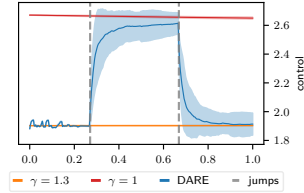


Figure 8: Mean and std dev of policy for DARE when the true value of γ jumps between 1.3 and 1.

334 4.6 Sample Efficiency

335 Finally, we show that DARE is more computationally efficient than the SOTA RL algorithms tested. In
 336 Table 1, we report the wall-clock training times of the offline phase each of the best performing RL
 337 methods and DARE in the LQG task. Clearly, DARE learns an optimal solution much faster than each
 338 of the RL methods. As the RL methods tested here are model-free, this is expected, as each algorithm
 339 must also learn from scratch the dynamics corresponding to the agent’s prior over the system. We
 340 leave the comparison of DARE to more model-based RL methods for future work.

Table 1: Average Wall-Clock Runtime Comparison, see Appendix J.

Algorithm	Wall-clock time in sec.	Device
A2C	451	NVIDIA A40
PPO	550	NVIDIA A40
SAC	1629	NVIDIA A40
DARE	22	Apple M1

341 5 Conclusions, Limitations, and Future Work

342 This paper presents DARE, a PINN-based methodology for solving OC problems in noisy and non-
 343 stationary environments. We demonstrate that when the time between actions is small, classical RL
 344 methods struggle in continuous-time control settings. Moreover, these methods are unstable with
 345 respect to action-repetition. In contrast, DARE does not require a fixed time discretization scheme and
 346 is shown to learn accurate solutions in stationary environments. Additionally, with the incorporation
 347 of a variance-adjusted objective which explicitly trades off exploration and exploitation, we show
 348 that DARE can efficiently adapt to non-stationary environments in real time.

349 A limitation of this work is the restriction of our experiments to simple control problems. Justification
 350 for our choice of problem settings is twofold. First, we sought to demonstrate that even simple
 351 continuous-time control problems prove difficult for classical RL methods without ad-hoc techniques.
 352 With these techniques, the RL methods considered were still unable to learn the true optimal control
 353 policy. Second, we sought OC problems for which classical solutions were available in order to verify
 354 that DARE indeed learns the true solution. We also wish to highlight that while the MPC environment
 355 appears simple, generating solution via classical methods is by no means trivial, as the system does
 356 not admit a dimensionality reducing ansatz. We acknowledge the need for further experimentation of
 357 DARE in physical systems and intend to pursue this in future work.

358 Also left for future work is a more unified methodology for system identification. In the LQG
 359 problem, the dynamics of the system could be learned quickly because of the problem’s simple linear
 360 structure. For more nonlinear problems, this estimation can be far more difficult. In future work, we
 361 intend to integrate DARE with more flexible methods for learning system dynamics.

References

- [1] A. Al-Arabi, A. Correia, G. Jardim, D. de Freitas Naiff, and Y. Saporito. Extensions of the deep galerkin method. *Applied Mathematics and Computation*, 430:127287, 2022.
- [2] F. Allgower, R. Findeisen, Z. K. Nagy, et al. Nonlinear model predictive control: From theory to application. *Journal-Chinese Institute Of Chemical Engineers*, 35(3):299–316, 2004.
- [3] A. S. Anand, J. E. Kveen, F. Abu-Dakka, E. I. Grøtli, and J. T. Gravdahl. Addressing sample efficiency and model-bias in model-based reinforcement learning. In *2022 21st IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 1–6. IEEE, 2022.
- [4] A. C. Aristotelous, E. C. Mitchell, and V. Maroulas. Adlgm: An efficient adaptive sampling deep learning galerkin method. *Journal of Computational Physics*, 477:111944, 2023.
- [5] A. Bachouch, C. Huré, N. Langrené, and H. Pham. Deep neural networks algorithms for stochastic control problems on finite horizon: numerical applications. *Methodology and Computing in Applied Probability*, 24(1):143–178, 2022.
- [6] M. Basei, X. Guo, A. Hu, and Y. Zhang. Logarithmic regret for episodic continuous-time linear-quadratic reinforcement learning over a finite-time horizon. *The Journal of Machine Learning Research*, 23(1):8015–8048, 2022.
- [7] A. Bensoussan. Perturbation methods in optimal control. (*No Title*), 1988.
- [8] T. Bhudisaksang and A. Cartea. Adaptive robust control in continuous time. *SIAM Journal on Control and Optimization*, 59(5):3912–3945, 2021.
- [9] K. Bieker, S. Peitz, S. L. Brunton, J. N. Kutz, and M. Dellnitz. Deep model predictive flow control with limited sensor data and online learning. *Theoretical and computational fluid dynamics*, 34:577–591, 2020.
- [10] T. Björk, M. H. Davis, and C. Landén. Optimal investment under partial information. *Mathematical Methods of Operations Research*, 71:371–399, 2010.
- [11] Á. Cartea, S. Jaimungal, and J. Penalva. *Algorithmic and high-frequency trading*. Cambridge University Press, 2015.
- [12] S. N. Cohen, C. Knochenhauer, and A. Merkel. Optimal adaptive control with separable drift uncertainty. *arXiv preprint arXiv:2309.07091*, 2023.
- [13] F. E. Dornier. Measuring progress in deep reinforcement learning sample efficiency. *arXiv preprint arXiv:2102.04881*, 2021.
- [14] K. Doya. Reinforcement learning in continuous time and space. *Neural computation*, 12(1): 219–245, 2000.
- [15] F. Drissi. Solvability of differential riccati equations and applications to algorithmic trading with signals. *Applied Mathematical Finance*, 29(6):457–493, 2022. doi: 10.1080/1350486X.2023.2241130. URL <https://doi.org/10.1080/1350486X.2023.2241130>.
- [16] G. Duran-Martin, A. Kara, and K. Murphy. Efficient online bayesian inference for neural bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 6002–6021. PMLR, 2022.
- [17] D. Evangelista and Y. Thamsten. Approximately optimal trade execution strategies under fast mean-reversion. *arXiv preprint arXiv:2307.07024*, 2023.
- [18] P. A. Forsyth and G. Labahn. Numerical methods for controlled hamilton-jacobi-bellman pdes in finance. *Journal of Computational Finance*, 11(2):1, 2007.
- [19] C. E. Garcia, D. M. Prett, and M. Morari. Model predictive control: Theory and practice—a survey. *Automatica*, 25(3):335–348, 1989.
- [20] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR, 2018.
- [21] J. Han, A. Jentzen, and W. E. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [22] M. Hoglund, E. Ferrucci, C. Hernandez, A. M. Gonzalez, C. Salvi, L. Sanchez-Betancourt, and Y. Zhang. A neural rde approach for continuous-time non-markovian stochastic control problems, 2023.

- 414 [23] C. Huré, H. Pham, A. Bachouch, and N. Langrené. Deep neural networks algorithms for
415 stochastic control problems on finite horizon: convergence analysis. *SIAM Journal on Numerical*
416 *Analysis*, 59(1):525–557, 2021.
- 417 [24] P. A. Iglesias and B. P. Ingalls. *Control theory and systems biology*. MIT press, 2010.
- 418 [25] Y. Jia and X. Y. Zhou. Policy gradient and actor-critic learning in continuous time and space:
419 Theory and algorithms. *The Journal of Machine Learning Research*, 23(1):12603–12652, 2022.
- 420 [26] Y. Jia and X. Y. Zhou. q-learning in continuous time. *Journal of Machine Learning Research*,
421 24(161):1–61, 2023.
- 422 [27] D. Jiang, J. Sirignano, and S. N. Cohen. Global convergence of deep galerkin and pinns methods
423 for solving partial differential equations. *arXiv preprint arXiv:2305.06000*, 2023.
- 424 [28] K. Kunisch and D. Walter. Semiglobal optimal feedback stabilization of autonomous systems via
425 deep neural network approximation. *ESAIM: Control, Optimisation and Calculus of Variations*,
426 27:16, 2021.
- 427 [29] I. Lenz, R. A. Knepper, and A. Saxena. Deepmpc: Learning deep latent features for model
428 predictive control. In *Robotics: Science and Systems*, volume 10, page 25. Rome, Italy, 2015.
- 429 [30] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.
430 Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- 431 [31] A. Mesbah. Stochastic model predictive control: An overview and perspectives for future
432 research. *IEEE Control Systems Magazine*, 36(6):30–44, 2016.
- 433 [32] H. N. Mhaskar. Neural networks for optimal approximation of smooth and analytic functions.
434 *Neural computation*, 8(1):164–177, 1996.
- 435 [33] P. K. Mishra, M. V. Gasparino, A. E. B. Velasquez, and G. Chowdhary. Deep model predictive
436 control, 2023.
- 437 [34] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller.
438 Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- 439 [35] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu.
440 Asynchronous methods for deep reinforcement learning. In *International conference on machine*
441 *learning*, pages 1928–1937. PMLR, 2016.
- 442 [36] A. Nagabandi, G. Kahn, R. S. Fearing, and S. Levine. Neural network dynamics for model-
443 based deep reinforcement learning with model-free fine-tuning. In *2018 IEEE international*
444 *conference on robotics and automation (ICRA)*, pages 7559–7566. IEEE, 2018.
- 445 [37] S. Niu, Y. Liu, J. Wang, and H. Song. A decade survey of transfer learning (2010–2020). *IEEE*
446 *Transactions on Artificial Intelligence*, 1(2):151–166, 2020.
- 447 [38] D. Onken, L. Nurbekyan, X. Li, S. W. Fung, S. Osher, and L. Ruthotto. A neural network
448 approach for high-dimensional optimal control applied to multiagent path finding. *IEEE*
449 *Transactions on Control Systems Technology*, 31(1):235–251, 2022.
- 450 [39] E. C. Ozcan, V. Giammarino, J. Queeney, and I. C. Paschalidis. A model-based approach for
451 improving reinforcement learning efficiency leveraging expert observations. *arXiv preprint*
452 *arXiv:2402.18836*, 2024.
- 453 [40] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on knowledge and*
454 *data engineering*, 22(10):1345–1359, 2009.
- 455 [41] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann. Stable-baselines3:
456 Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22
457 (268):1–8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>.
- 458 [42] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep
459 learning framework for solving forward and inverse problems involving nonlinear partial
460 differential equations. *Journal of Computational physics*, 378:686–707, 2019.
- 461 [43] T. Salzmann, E. Kaufmann, J. Arrizabalaga, M. Pavone, D. Scaramuzza, and M. Ryll. Real-time
462 neural MPC: Deep learning model predictive control for quadrotors and agile robotic platforms.
463 *IEEE Robotics and Automation Letters*, 8(4):2397–2404, apr 2023. doi: 10.1109/lra.2023.
464 3246839. URL <https://doi.org/10.1109/2F1ra.2023.3246839>.

- 465 [44] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
466 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 467 [45] J. Sirignano and K. Spiliopoulos. Dgm: A deep learning algorithm for solving partial differential
468 equations. *Journal of computational physics*, 375:1339–1364, 2018.
- 469 [46] N. Srinivas, A. Krause, S. M. Kakade, and M. Seeger. Gaussian process optimization in the
470 bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*, 2009.
- 471 [47] P. M. Suder, J. Xu, and D. B. Dunson. Bayesian transfer learning. *arXiv preprint*
472 *arXiv:2312.13484*, 2023.
- 473 [48] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 474 [49] C. Tallec, L. Blier, and Y. Ollivier. Making deep q-learning methods robust to time discretization.
475 In *International Conference on Machine Learning*, pages 6096–6104. PMLR, 2019.
- 476 [50] C. Tan, F. Sun, T. Kong, W. Zhang, C. Yang, and C. Liu. A survey on deep transfer learning. In
477 *Artificial Neural Networks and Machine Learning–ICANN 2018: 27th International Conference*
478 *on Artificial Neural Networks, Rhodes, Greece, October 4-7, 2018, Proceedings, Part III 27*,
479 pages 270–279. Springer, 2018.
- 480 [51] E. Todorov and W. Li. A generalized iterative lqg method for locally-optimal feedback control
481 of constrained nonlinear stochastic systems. In *Proceedings of the 2005, American Control*
482 *Conference, 2005.*, pages 300–306. IEEE, 2005.
- 483 [52] R. van der Meer, C. W. Oosterlee, and A. Borovykh. Optimally weighted loss functions for
484 solving pdes with neural networks. *Journal of Computational and Applied Mathematics*, 405:
485 113887, 2022.
- 486 [53] H. Wang, T. Zariphopoulou, and X. Y. Zhou. Reinforcement learning in continuous time and
487 space: A stochastic control approach. *The Journal of Machine Learning Research*, 21(1):
488 8145–8178, 2020.
- 489 [54] C. Yildiz, M. Heinson, and H. Lähdesmäki. Continuous-time model-based reinforcement
490 learning. In *International Conference on Machine Learning*, pages 12009–12018. PMLR, 2021.
- 491 [55] J. Yong and X. Y. Zhou. *Stochastic controls: Hamiltonian systems and HJB equations*, vol-
492 ume 43. Springer Science & Business Media, 1999.
- 493 [56] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He. A comprehensive
494 survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020.

495 **A Related Work**

496 **Deep Learning Methods for Control.** Deep learning is extensively used to solve HJB equations
 497 because of its flexibility and scalability to high dimensions [23, 5, 38, 28]. Among earlier examples,
 498 [21, 45] provide two contrasting approaches. The former proposes the Deep Galerkin Method, which
 499 uses Monte Carlo integration to minimize a variational form of the HJB in a mesh-free manner, while
 500 the latter proposes the Deep BSDE method, which reformulates the PDE as a backward stochastic
 501 differential equation and approximates the gradient of the solution with a neural network. Global
 502 convergence of DGM was recently established in [27]. There are numerous extensions to their
 503 method, including adaptive Monte Carlo sampling in [4], augmented loss function for non-parametric
 504 running penalties and drifts in [1], and optimally weighted loss objectives in [52].

505 **Model Predictive Control.** Classical methods in MPC are the foundation of many online control
 506 optimization methods in both the deterministic and stochastic settings [2, 31]. Recently, deep learning
 507 was integrated with MPC, with applications in controlling uncertain nonlinear systems such as
 508 unsteady fluid flow and high-performance autonomous systems [29, 33, 9, 43, 36]. These approaches
 509 leverage neural networks to enhance dynamic modeling capacity and real-world control performance.

510 **Continuous-Time Reinforcement Learning.** Methods in deep RL are highly effective in several
 511 complex decision-making problems [34, 30, 44]. Continuous-time environments pose significant
 512 challenges to RL methods [49, 54]. In particular many RL methods optimize incorrect objectives
 513 [see 25] when environments are noisy, e.g., temporal difference (TD) learning [14]. Recently, [53]
 514 study the exploration-exploitation trade-off in stochastic and continuous-time RL, and prove that the
 515 optimal exploration policy is Gaussian in a Linear-Quadratic setting. Subsequent work [25, 26, 22, 6]
 516 extend RL methods to stochastic and continuous-time environments.

517 **B Algorithms**

518

Algorithm 1 : DARE - Offline Phase

Inputs:

- Initial OC pair $\widehat{\mathbf{p}}_0 = (\widehat{h}_0, \widehat{f}_0)$
- Weights θ_0, ψ_0
- Area of integration $K \subset \mathbb{R}$
- Number of training iterations $N \in \mathbb{N}$
- Initial state X_0

$$\widehat{u}_{X_0} \leftarrow \text{ILQG}(X_0, \widehat{\mathbf{p}}_0)$$

$$V^{\theta_0} \leftarrow \widehat{g}_0 + \mathfrak{X}^{\theta_0}$$

$$u^{\psi_0} \leftarrow \widehat{u}_{X_0} + \mathfrak{X}^{\psi_0}$$

for $n = 1, \dots, N$ **do**

$\mathcal{D} \leftarrow$ *Batch of uniform samples of* $[0, T) \times K$

$$\ell \leftarrow \frac{1}{|\mathcal{D}|} \sum_{(t,x) \in \mathcal{D}} \mathcal{L}(t, x, V^{\theta_0}, u^{\psi_0}, \widehat{\mathbf{p}}_0)$$

$$\theta_n, \psi_n \leftarrow \text{ADAM}(\theta_{n-1}, \psi_{n-1}, \text{loss} = \ell)$$

end for

Algorithm 2 : DARE - Online Phase

Input:

- OC pair $\widehat{\mathbf{p}}_0 = (\widehat{h}_0, \widehat{f}_0)$
- Initial weights θ, ψ
- Area of integration $K \subset \mathbb{R}$
- Number of ADAM updates in offline step $N_{off} \in \mathbb{N}$
- Number of ADAM updates per online step $N_{on} \in \mathbb{N}$
- Time discretization $\mathcal{T} \subset (0, T)$
- Initial State X_0

$$V^{\theta_0}, u^{\psi_0} \leftarrow \text{DARE}(\widehat{\mathbf{p}}_0, \theta, \psi, K, N_{off}, X_0)$$

$$t_{prev} \leftarrow 0$$

$$\widehat{\mathbf{p}}_{prev} \leftarrow \widehat{\mathbf{p}}_0$$

for $t \in \mathcal{T}$ **do**

$$X_t \leftarrow \text{System}(t, X_{t_{prev}}, u^{\psi_{t_{prev}}}(t, X_{t_{prev}}))$$

$$\widehat{\mathbf{p}}_t \leftarrow \text{Approx}(\widehat{\mathbf{p}}_{prev}, \widehat{\mathbf{p}}_{prev}(X_t, u_{t_{prev}}^{\psi}(t, X_t))) + \epsilon)$$

$$V^{\theta_t}, u^{\psi_t} \leftarrow$$

$$\text{DARE}(\widehat{\mathbf{p}}_t, \theta_{t_{prev}}, \psi_{t_{prev}}, t, N_{on})$$

$$t_{prev} \leftarrow t$$

end for

519 **C Experimental Parameters**

520 Default parameters for the LQG and MPC problems are reported here.

Table 2: Default parameter values for the LQG problem.

PARAM.	b	c	σ	ϕ	α	x_0	T
VALUE	-5	1	1	1	0.3	10	1

Table 3: Default parameters values for the MPC problem.

Param.	b	c	σ	ϕ	φ	α	x_0	T
Value	0	1	1	0.15	0.1	0.05	100	1

521 **D ILQG**

522 We provide a brief overview of the ILQG method from [51] that we use to initialize the control policy.
 523 Let the system X_t evolve as:

$$dX_t = h(X_t, u_t)dt + \Sigma(X_t, u_t)dW_t$$

524 and let the performance criterion be

$$J(t, x; u) = \mathbb{E} \left[g(X_T) + \int_t^T f(\tau, X_\tau, u_\tau) d\tau \right].$$

525 In this section, we assume that the agent seeks to *minimize* $J(t, x; u)$. Let \bar{u}_t be a random *open-loop*
 526 control policy, and consider

$$d\bar{X}_t = f(\bar{X}_t, \bar{u}_t).$$

527 Next, we linearize the original system around \bar{X}_t, \bar{u}_t and discretize time $k = \{0, \dots, K\}$ with
 528 $\Delta t = \frac{T}{K-1}$ and $t_k = k\delta t$.

529 Define the discrepancies $\delta X_t = X_t - \bar{X}_t, \delta u_t = u_t - \bar{u}_t$, which evolve (approximately) as

$$\begin{aligned} \delta X_{k+1} &= A_k \delta X_k + B_k \delta u_k + C_k (\delta u_k) \xi_k \\ C_k &= c_{1,k} + C_{1,k} \delta u_k + \dots + C_{d,d_u} \\ \text{cost}_k &= q_k + \delta X_k^\top \mathbf{q}_k + \frac{1}{2} \delta X_k^\top Q_k \delta X_k \\ &\quad + \delta u_k^\top \mathbf{r}_k + \frac{1}{2} \delta u_k^\top R_k \delta u_k + \delta u_k^\top P_k \delta X_k, \end{aligned}$$

530 where $\delta X_0 = 0, \xi_k \sim N(0, I_{d_x})$,

$$\begin{aligned} A_k &= I_{d_x} + \Delta t h_x & \mathbf{q}_k &= \Delta t f_x \\ B_k &= \Delta t h_u & Q_k &= \Delta t f_{xx} \\ c_{i,k} &= \sqrt{\Delta t} \Sigma^i & \mathbf{r}_k &= \Delta t f_u \\ C_{i,k} &= \sqrt{\Delta t} \Sigma_u^i & R_k &= \Delta t f_{uu} \\ q_k &= \Delta t f & P_k &= \Delta t f_{ux}, \end{aligned}$$

531 and $q_K = g, \mathbf{q}_K = g_x$, and $Q_K = g_{xx}$.

532 Above, all functions are evaluated at \bar{X}_k, \bar{u}_k , and Σ^i denotes the i -th row of Σ . It is shown in [51]
 533 that the optimal control to the linearized system δu^* is affine, with

$$\delta u^*(\delta X) = l_k + L_k \delta X. \tag{16}$$

534 When δu takes the form (16), the value function is quadratic and we write

$$V_k(\delta X) = s_k + \delta X_k^\top \mathbf{s}_k + \frac{1}{2} \delta X_k^\top S_k \delta X_k.$$

535 One can obtain an explicit representation of S_k, \mathbf{s}_k, s_k by first defining

$$\begin{aligned} \mathbf{g}_k &= \mathbf{r}_k + B_k^\top \mathbf{s}_k + \sum_i C_{i,k}^\top S_{k+1} c_{i,k} \\ G_k &= P_k + B_k^\top S_{k+1} A_k \\ H_k &= R_k + B_k^\top B_k + \sum_i C_{i,k}^\top S_{k+1} C_{i,k}, \end{aligned}$$

536 which leads to the following equalities

$$\begin{aligned} S_k &= Q_k + A_k^\top S_{k+1} A_k - L_k^\top H_k L_k + L_k^\top G_k + G_k^\top L_k \\ \mathbf{s}_k &= \mathbf{q}_k + A_k^\top \mathbf{s}_{k+1} + L_k^\top H_k l_k + L_k^\top \mathbf{g}_k + G_k^\top l_k \\ s_k &= q_k + s_{k+1} + \frac{1}{2} \sum_i c_{i,k} S_{k+1} c_{i,k} + \frac{1}{2} l_k^\top H_k l_k + l_k^\top \mathbf{g}_k, \end{aligned}$$

537 where $S_K = Q_K, \mathbf{s}_K = \mathbf{q}_K, s_K = q_K$. Consequently, we obtain

$$\begin{aligned} l_k &= -H_k^{-1} \mathbf{g}_k \\ L_k &= -H_k^{-1} G_k. \end{aligned}$$

538 When f or g are not convex, H may have negative eigenvalues. This generally causes numerical
539 issues due to the the minimization problem being unbounded. In this case, we use the Levenberg-
540 Marquardt method to achieve an approximate inverse, by forcing all negative eigenvalues of H to be
541 equal to some $\lambda > 0$.

542 E Transfer Learning in Neural Adaptive Control

543 One way of interpreting the ability of DARE to adapt to unseen environments in real time is through
544 transfer Learning (TL). TL encompasses methods in which knowledge acquired from an initial source
545 task is used to improve performance on a related target task; see [40, 56, 37, 47, 50] for an overview
546 of transfer learning. One can study the efficiency of DARE in the online phase from the perspective of
547 TL because its performance hinges on successive transferring of knowledge (parameters) between
548 DNNs corresponding to the solutions to “similar” OC problems; see Figure 1. In this section, we
549 provide a theoretical justification for our method. More precisely, we analyze the smoothness of OC
550 problems with respect to the OC pair describing the environment and the resulting smoothness of
551 DNN parameters.

552 To provide a theoretical foundation to this claim, we use the tools of regular perturbation in OC
553 and a notion of continuity of the DARE network parameters. Later, Section 4.2 explores specific
554 examples and quantifies empirically the improvement achieved from TL in the online phase of DARE.
555 In particular, we use the number of iterations required to attain, on average, a prespecified loss in the
556 target task to measure the *strength of transfer*.

557 To streamline our analysis, assume $d_X = d_u = 1$ and consider an agent who receives observations of
558 the OC pair and updates their estimate $\hat{\mathbf{p}}_t$, accordingly.² In practice, between two sufficiently close
559 observation times $r, s \in [0, T]$ with $r < s$, we assume that the estimate $\hat{\mathbf{p}}_r$ at time r remains close
560 to the estimate $\hat{\mathbf{p}}_s$ at time s . Hence, we write $\hat{\mathbf{p}}_s$ as a perturbation of $\hat{\mathbf{p}}_r$. This is formalized in the
561 following assumption.

562 **Assumption E.1.** For any ϵ , there are suitable perturbation functions p^f and p^h such that

$$\hat{f}_s = \hat{f}_r + \epsilon p^f \quad \text{and} \quad \hat{h}_r = \hat{h}_r + \epsilon p^h. \quad (17)$$

²It is straightforward to generalize to multi-dimensional setups.

563 Fix $t \in [0, T]$ and consider the value and control functions associated to $\widehat{\mathbf{p}}_r$ and $\widehat{\mathbf{p}}_s$ on $[t, T]$. That is,
 564 for $\rho \in \{r, s\}$, let

$$V^\rho(t, x) = \sup_u \mathbb{E} \left[\widehat{g}(X_T^\rho) + \int_t^T \widehat{f}_\rho(X_\tau^\rho, u_\tau) d\tau \mid X_t^\rho = x \right], \quad (18)$$

565 where

$$dX_\tau^\rho = \widehat{h}_\rho(X_\tau^\rho, u_\tau) d\tau + \widehat{\Sigma} dW_\tau.$$

566 Theorem E.2 first shows that small perturbations in the OC pair lead to small perturbations in the
 567 optimal policy and the value function. Next, the result examines the continuity of the parameters of
 568 the DNNs approximating the value and control functions. This continuity is considered with respect
 569 to the function space in which the functions are defined. Intuitively, when a DNN is trained to a
 570 particular function, one expects that marginal changes to this function will result in marginal changes
 571 to the network parameters. Providing such a result in a general setting poses an intricate challenge.
 572 Thus, we simplify the setting by reducing the class of DNNs to that of single-layer perceptrons.
 573 However, our empirical findings suggest that it generalizes to more general cases.

574 **Theorem E.2.** *Suppose that $\widehat{f}^r, \widehat{f}^s, \widehat{h}^r, \widehat{h}^s, p^f, p^h \in C_b^{1,2}([0, T]; K)$ and ϵ is defined as in (17).³
 575 Moreover, assume that solutions $(V^r, u^{r,*})$ and $(V^s, u^{s,*})$ to (18) exist and are unique. For a value
 576 of ϵ which is sufficiently small, then there exists $L > 0$ such that for any $\gamma > 0$ and single-layer
 577 perceptron approximations $(V^{\theta^r}, u^{\psi^r})$ and $(V^{\theta^s}, u^{\psi^s})$ of (V^r, u^r) and (V^s, u^s) , respectively, with
 578 precision γ , such that*

$$\|\theta^r - \theta^s\| + \|\psi^r - \psi^s\| \leq L\epsilon^2. \quad (19)$$

579 The proof of Theorem E.2 is given below. Although the above result justifies the performance of
 580 DARE for two consecutive policy updates with two fixed OC pairs, the results extend to the case
 581 of a dynamic estimate of the environment. That is, suppose $(h, f) = (h_t, f_t)$ evolves throughout
 582 $t \in [0, T]$. Then, if (h_t, f_t) changes smoothly, we expect the DNNs parameterizing the corresponding
 583 solutions to vary smoothly. Finally, our numerical results indicate that DARE also adapts efficiently to
 584 large and abrupt changes in the environment.

585 E.1 Definitions

586 First, we introduce the notation used throughout the section.

587 **Definition E.3** (Single-Layer Perceptron (SLP)). Denote $d_i, d_h, d_o \in \mathbb{N}, \sigma : \mathbb{R} \rightarrow \mathbb{R}$. A *single-layer*
 588 *perceptron* is defined as

$$F : \begin{array}{l} \mathbb{R}^{d_i} \longrightarrow \mathbb{R}^{d_o} \\ x \longmapsto \sum_{i=1}^{d_h} (C^\top)_i \phi \bullet (A_i x + b_i) \end{array}$$

589 with $A_i \in \mathbb{R}^{d_i}, b_i \in \mathbb{R}, (C^\top)_i \in \mathbb{R}^{d_o}$ for $i \in \{1, \dots, d_h\}$ and \bullet denotes the component-wise
 590 application. We denote with $\theta := (A, b, C) \in \mathbb{R}^d$ the parameters of this SLP, with $d = d_i d_h + d_h +$
 591 $d_o d_h$, and F_θ is an SLP with parameter θ .

592 E.2 Proof of Theorem E.2

593 We split the proof of Theorem E.2 into two results. Proposition E.4 shows that perturbations in the
 594 environment lead to perturbations of similar scale in the value function and the optimal policy of OC
 595 problems. Next, Proposition E.5 shows that for small perturbations of the value function and optimal
 596 policy, the parameters of the networks used to approximate these functions are continuous.

597 **Proposition E.4.** *There is a constant C such that*

$$|V^r(t, x) - V^s(t, x)| \leq C \epsilon^2, \quad (20)$$

$$|u^{r,*} - u^{s,*}| \leq C \epsilon. \quad (21)$$

³ $C_b^{1,2}([0, T]; K)$ denotes the set of functions defined on $[0, T] \times K$ with continuous first derivative in t and continuous and bounded second derivatives in x .

598 **Proof** First, note that the assumption of Theorem E.2 ensure that the functional J is well defined.
 599 Observe that the OC problem V^r is a perturbation of V^s but also V^s is a perturbation of V^r . In
 600 particular, first write

$$J^r(t, x, u) = \mathbb{E} \left[\widehat{g}(X_T^s) + \int_t^T \widehat{f}_s(X_\tau^s, u_\tau) d\tau \mid X_t^s = x \right] \quad (22)$$

$$J^r(t, x, u) = \mathbb{E} \left[\widehat{g}(X_T^r) + \int_t^T \widehat{f}_s(X_\tau^r, u_\tau) d\tau \mid X_t^r = x \right]. \quad (23)$$

601 Next, use Theorem 2.1 of Chapter III and Theorem 2.1 of Chapter IV in [7] to write

$$\begin{cases} |V^r(t, x) - J^r(t, x, u^{s,*})| & \leq C_0 \epsilon^2, \\ |V^s(t, x) - J^s(t, x, u^{r,*})| & \leq C_1 \epsilon^2, \end{cases}$$

602 for suitable constants C_0 and C_1 . Finally, use the asymptotic expansions (Section 2.3, Chapter III, in
 603 [7]) to write

$$\begin{cases} |V^r(t, x) - J^s(t, x, u^{s,*})| & \leq L_0 \epsilon^2, \\ |V^s(t, x) - J^r(t, x, u^{r,*})| & \leq L_1 \epsilon^2, \end{cases}$$

604 for suitable constants L_0 and L_1 . □

605 **Proposition E.5.** *Let $K \subseteq \mathbb{R}^{d_i}$ be compact; $f \in C_b^2(K; \mathbb{R})$. There exists $\delta, L > 0$ such that for every
 606 $f' : K \rightarrow \mathbb{R}^{d_o}$ with $\|f - f'\| < \delta$ and every $\gamma > 0$, there exists parameters $\theta, \theta' \in \mathbb{R}^d$ such that*

$$\|F_\theta - f\|_{C_b^2(K; \mathbb{R})} \leq \gamma, \quad (24)$$

$$\|F_{\theta'} - f'\|_{C_b^2(K; \mathbb{R})} \leq \gamma, \quad (25)$$

607 and

$$\|\theta' - \theta\| < L \|f - f'\|_{C_b^2(K; \mathbb{R})}, \quad (26)$$

608 where F_θ is a single-layer perceptron with ReLU activation of width d_h .

609 **Proof** Without loss of generality, assume that K includes an open set around 0, i.e. there exists
 610 $0 < \delta < \epsilon$ s.t. $\mathcal{D}^\delta - f \subseteq C_b^2(K; \mathbb{R})$, where

$$\mathcal{D}^\delta := \left\{ f' \in C_b^2(K; \mathbb{R}) \mid \|f - f'\|_{C_b^2(K; \mathbb{R})} < \delta \right\}.$$

611 Fix an arbitrary $\gamma > 0$. According to [32], Theorem 2.1, we can find a hidden dimension $d_h \in \mathbb{N}$, a
 612 matrix $A \in \mathbb{R}^{d_h \times d_i}$, a vector $b \in \mathbb{R}^{d_h}$, and a continuous linear functional $\mathcal{C} : C_b^2(K; \mathbb{R}) \rightarrow \mathbb{R}^{d_h}$ such
 613 that

$$\|f' - F_{(A,b,\mathcal{C}(f'))}\|_p \leq \gamma, \quad f \in \mathcal{D}^\delta.$$

614 Since for $f' \in \mathcal{D}^\delta$ holds

$$F_{(A,b,\mathcal{C}(f'))} = F_{(A,b,\mathcal{C}(f'-f))} + F_{(A,b,\mathcal{C}(f))},$$

615 due to linearity of \mathcal{C} and Definition E.3, we have

$$\|\theta' - \theta\| = \|\mathcal{C}(f' - f)\| \leq L \|f - f'\|_{C_b^2(K; \mathbb{R})}$$

616 where we used the fact that the operator norm $\|\mathcal{C}\|_T$ of a continuous operator is finite. We conclude
 617 $\|\theta' - \theta\| < \epsilon$ and finish the proof. □

618 F Performance

619 We report training performance of all the methods tested in Section 4.1 in Table 4. All tests have
 620 been conducted on a standard MacBook Pro M1 and a single NVIDIA A40 for training each
 621 of the RL algorithms. We make our code public in the following (anonymized) repo: <https://anonymous.4open.science/r/dare-7136/README.md>
 622

Table 4: Run time of different algorithms in the experiments.

Section	Test	Algorithm	Run time in seconds per 1000 iteration
4.2 Training performance	Offline LQG	DGM	4.26
		MLP	1.95
		DARE	2.24
	Offline MPC	DGM	6.02
		MLP	2.28
		DARE	2.34
4.4 LQG	Online phase	DARE	7.47
4.5 MPC	Exploration-Exploitation Non-stationary	DARE	2.46
		DARE	2.52
4.1 High-dimensional	Offline phase	DGM	6.01
		MLP	2.41
		DARE	4.01
	Online phase	DGM	6.64
		MLP	3.26
		DARE	3.92

623 G Filtering mathematics

624 G.1 Perfect knowledge of the drift

625 This section solves the OC problem (14) when the agent fixed the value of the drift and does not
626 update their belief throughout the time window.

627 When the drift is known and fixed, the OC problem (14) can be solved with standard methods [55],
628 and is

$$u^* = \frac{c}{2\phi} (2A(t)x + B(t) + 1), \quad (27)$$

629 where A and B solve the ODE system

$$\begin{cases} -A'(t) = \frac{cA(t)^2}{2\phi} \\ -B'(t) = 2\mu A(t) + \frac{c^2 A(t)(B(t)+1)}{\phi} \end{cases}. \quad (28)$$

630 G.2 Bayesian filtering of the Gaussian drift

631 This section solves the OC problem when the agent uses a Gaussian prior to continuously update
632 their estimation of the drift throughout the time window of the OC problem.

633 Consider the control problem in (14). When the agent uses a Gaussian prior $\mathcal{N}(b_t, \Pi_0)$ for μ then it
634 can be shown that the dynamics of x can be written

$$dx_t = \beta_t dt + c u_t dt + \sigma d\widehat{W}_t$$

in a different filtration in which \widehat{W} is a Gaussian process. $\beta_t = \mathbb{E}[\mu | \mathcal{F}_t]$ is the best estimate of μ at time t and can be obtained analytically as

$$\beta_t = -\frac{\Pi(t)}{\sigma} \left(x_0 - \frac{\sigma b_0}{\Pi_0} - x_t + q_t \right)$$

635 where $\Pi(t) = (\Pi_0^{-1} + \frac{t}{\sigma})^{-1}$ and $q_t = \int_0^t c u_t dt$.

636 Using the learning dynamics above to solve the control problem (see [15] for details) gives the
637 optimal control \tilde{u}^* given by

$$\tilde{u}^* = \frac{c}{2\phi} (2A(t) + B(t))x + (2C(t) + B(t))q \quad (29)$$

$$+ (1 + D(t) + E(t)), \quad (30)$$

where A, B, C solve the Riccati equation in

$$\begin{aligned} P(t) &= \begin{pmatrix} A(t) & \frac{1}{2}B(t) \\ \frac{1}{2}B(t)^\top & C(t) \end{pmatrix} \\ 0 &= P'(t) + Y(t)^\top P(t) + P(t) Y(t) + P(t) U P(t), \\ Y(t) &= \begin{pmatrix} \frac{\Pi(t)}{\sigma} & \frac{\Pi(t)}{\sigma} \\ \frac{\sigma}{\sigma} & \frac{\sigma}{\sigma} \end{pmatrix}, \quad U = \begin{pmatrix} \frac{c^2}{\phi} & \frac{c^2}{\phi} \\ \frac{c^2}{\phi} & \frac{c^2}{\phi} \end{pmatrix}, \\ P(T) &= \begin{pmatrix} -\alpha & 0 \\ 0 & 0 \end{pmatrix}, \end{aligned}$$

and D and E solve the ODE system

$$\begin{cases} 0 = D'(t) + 2\bar{\Pi} \Pi(t) A(t) - \Pi(t) (1 + D(t)) \\ \quad + \frac{c^2}{4\phi} (2A(t) + B(t))^2 \\ 0 = E'(t) + \bar{\Pi} \Pi^\top B(t) + \Pi^\top (1 + D(t)) \\ \quad + \frac{c^2}{4\phi} (2C(t) + B(t))^2, \end{cases}$$

638 with terminal conditions $D(T) = E(T) = 0$.

639 H Algorithmic trading in high dimension

640 We motivate the multidimensional setup in our experiments of Section 4.1. Consider the case of
641 the trading desk of a large bank that must execute a number $d \in \mathbb{N}^*$ of large transactions in d
642 correlated financial assets throughout a trading window $[0, T]$. The trading desk must minimize their
643 trading costs while minimizing the risk of their positions. Throughout this section, we consider a
644 filtered probability space $(\Omega, \mathcal{F}, \mathbb{P}; \mathbb{F} = (\mathcal{F}_t)_{t \in [0, T]})$, with $T > 0$, satisfying the usual conditions and
645 supporting all the processes we introduce.

646 Let $Q_0 \in \mathbb{R}^d$ represent the transaction sizes in every asset. The inventory of the agent is mod-
647 eled by $(Q_t)_{t \in [0, T]} = (Q_t^1, \dots, Q_t^d)_{t \in [0, T]}^\top$ and it evolves with the trading speed $(u_t)_{t \in [0, T]} =$
648 $(u_t^1, \dots, u_t^d)_{t \in [0, T]}^\top$ in each asset:⁴

$$dQ_t = u_t dt.$$

649 The prices $(S_t)_{t \in [0, T]} = (S_t^1, \dots, S_t^d)_{t \in [0, T]}^\top$ of the d assets are modeled as correlated Brownians
650 with dynamics

$$dS_t = \tilde{\Sigma} dW_t,$$

651 where $W = (W^1, \dots, W^d)$ is a d -dimensional standard Brownian motion and $S_0 \in \mathbb{R}^d$ is known.

652 The matrix $\tilde{\Sigma} \in \mathcal{M}_d(\mathbb{R})$ measures the correlation of the prices and we define the covariance matrix
653 $\Sigma = \tilde{\Sigma} \tilde{\Sigma}^\top \in \mathcal{S}_d^{++}(\mathbb{R})$.⁵

654 Trading activity of the agent generates transaction costs, driven by some function of the trading speed
655 $f(u_t)$ so the cash from their trading activity evolves as

$$dX_t = -u_t^\top S_t dt - f(u_t) dt, \quad X_0 = 0.$$

656 The agent maximizes the exponential utility of their terminal wealth so their objective is

$$V(t, x, q, s) = \sup_v \mathbb{E} \left[- \exp \left(- \gamma (Q_T^\top S_T - Q_T^\top \Gamma Q_T) \right. \right. \quad (31)$$

$$\left. \left. - \int_t^T u_s^\top S_s ds - \int_t^T f(u_s) ds \right) \right], \quad (32)$$

⁴The superscript $^\top$ is the transpose operator.

⁵ $\mathcal{M}_d(\mathbb{R}) := \mathcal{M}_{d,d}(\mathbb{R})$ is the set of $d \times d$ real square matrices, $\mathcal{S}_d(\mathbb{R})$ is the set of real symmetric $d \times d$ matrices, and $\mathcal{S}_d^{++}(\mathbb{R})$ is the set of positive matrices.

657 for values $Q_t = q$, $X_t = x$, and $S_t = s$ at time t .

658 The dynamic programming principle holds and the HJB equation associated with the problem

$$0 = \partial_t V + \frac{1}{2} \text{Tr} (\Sigma D_{SS}^2 V) + \sup_{u \in \mathbb{R}^d} (-(u^\top s + f(u)) \partial_x V + v^\top \nabla_q V), \quad (33)$$

659 with terminal condition

$$V(T, x, q, s) = -\exp(-\gamma (q^\top s - q^\top \Gamma q)). \quad (34)$$

660 In the experiment of Section 4.1, we solve the HJB (33)-(34) using DARE to obtain the optimal policy
661 of the trading agent.

662 When all the parameters of the problem are known and fixed, i.e., the agent does not adapt to new
663 information, the problem described above admits an analytical solution which we use to study the
664 performance of DARE.

665 To solve the problem semi-analytically, the function f must be a quadratic form, that is, there is some
666 $\eta \in S_d^{++}(\mathbb{R}^{d \times d})$ with

$$f(u) = u^\top \eta u. \quad (35)$$

667 We follow the standard steps in linear-exponential quadratic Gaussian (LEQG) control and we propose
668 the following form for the value function

$$V(t, x, q, s) = -\exp(-\gamma (x + q^\top S + Q^\top A(t)q + B(t)^\top q + C(t))),$$

669 and straightforward calculations find that the problem reduces to solving the following ODE system

$$\begin{cases} A'(t) = \frac{\gamma}{2} \Sigma - A(t) \eta^{-1} A(t) \\ B'(t) = -A(t) \eta^{-1} B(t) \\ C'(t) = -\frac{1}{4} B(t)^\top \eta^{-1} B(t), \end{cases} \quad (36)$$

670 with terminal conditions

$$A(T) = -\Gamma, \quad B(T) = C(T) = 0. \quad (37)$$

Clearly, the solutions for B and C are $B = C = 0$. To obtain a solution, we use the change of variables

$$a(t) = \eta^{-\frac{1}{2}} A(t) \eta^{-\frac{1}{2}} \quad \forall t \in [0, T],$$

671 so the problem reduces to the following terminal value problem

$$\begin{cases} a'(t) = \hat{A}^2 - a(t)^2 \\ a(T) = -C, \end{cases} \quad (38)$$

where

$$\hat{A} = \sqrt{\frac{\gamma}{2}} \left(\eta^{-\frac{1}{2}} \Sigma \eta^{-\frac{1}{2}} \right)^{\frac{1}{2}} \in S_d^{++}(\mathbb{R}),$$

and

$$C = \eta^{-\frac{1}{2}} \Gamma \eta^{-\frac{1}{2}} \in S_d^+(\mathbb{R}).$$

672 We solve (38) in the next result.

673 **Proposition H.1.** Define $\xi : [0, T] \rightarrow S_d(\mathbb{R})$

$$\begin{aligned} \xi(t) = & -\frac{\hat{A}^{-1}}{2} \left(I - e^{-2\hat{A}(T-t)} \right) \\ & - e^{-\hat{A}(T-t)} \left(C + \hat{A} \right)^{-1} e^{-\hat{A}(T-t)} \end{aligned} \quad (39)$$

674 as the unique solution to the ODE system

$$\begin{cases} \xi'(t) = \hat{A}\xi(t) + \xi(t)\hat{A} + I_d \\ \xi(T) = - (C + \hat{A})^{-1}. \end{cases} \quad (40)$$

Then $\forall t \in [0, T]$, $\xi(t)$ is invertible and

$$a : t \in [0, T] \rightarrow \hat{A} + \xi(t)^{-1} \in \mathcal{S}_d(\mathbb{R})$$

675 is the unique solution of (38).

676 Thus, the value function, which we use as the oracle in Section 4.1 is given by

$$V(t, x, q, s) = - \exp(-\gamma(x + q^T S + Q^T A(t)q + B(t)^T q + C(t))),$$

677 where

$$A(t) = \eta^{\frac{1}{2}} \left(\hat{A} - \left\{ \frac{\hat{A}^{-1}}{2} (I - e^{-2\hat{A}(T-t)}) + e^{-\hat{A}(T-t)} (C + \hat{A})^{-1} e^{-\hat{A}(T-t)} \right\}^{-1} \right) \eta^{\frac{1}{2}}.$$

678 Finally, Figure H shows the true value function and the solution learned by DARE for a set of model
679 parameters in dimension 5, and Figure 9 shows the associated training loss.

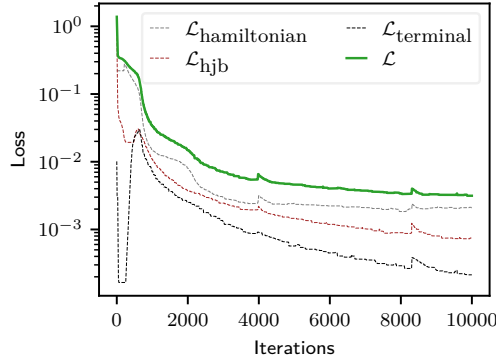


Figure 9: Training loss (10) of DARE for the multidimensional OC problem (31). The parameter values are $\Gamma = 10^{-2} \times I_5$,

$$\eta = 10^{-3} \times \begin{pmatrix} 25.13 & 10.41 & 11.67 & 13.75 & 22.21 \\ 10.41 & 6.42 & 7.68 & 9.12 & 12.4 \\ 11.67 & 7.68 & 12.95 & 11.2 & 18.71 \\ 13.75 & 9.12 & 11.2 & 17.04 & 17.25 \\ 22.21 & 12.4 & 18.71 & 17.25 & 29.02 \end{pmatrix}, \quad \text{and} \quad \Sigma = \begin{pmatrix} 2.83 & 2.02 & 2.53 & 1.91 & 1.59 \\ 2.02 & 1.96 & 2.04 & 1.28 & 1.31 \\ 2.53 & 2.04 & 2.85 & 2.04 & 1.32 \\ 1.91 & 1.28 & 2.04 & 1.76 & 0.96 \\ 1.59 & 1.31 & 1.32 & 0.96 & 1.06 \end{pmatrix}.$$

680 I Gaussian Process mathematics

Formally, a GP is a random function $f : \mathcal{X} \mapsto \mathbb{R}$, such that, for any finite set of points $\mathbf{X}_* \subseteq \mathcal{X}$, the random vector $\mathbf{f}_* = \{f(x)\}_{x \in \mathbf{X}_*}$ follows a multivariate Gaussian distribution. The shape of the function f is determined by a finite set of (training) observations $\mathbf{y} = \{y_i\}_{i \in \{1, \dots, n\}}$ collected at the (training) observation points $\mathbf{X} = \{\mathbf{x}_i\}_{i \in \{1, \dots, n\}}$, where $y_i = f(\mathbf{x}_i) + \epsilon_i$ is subject to i.i.d.

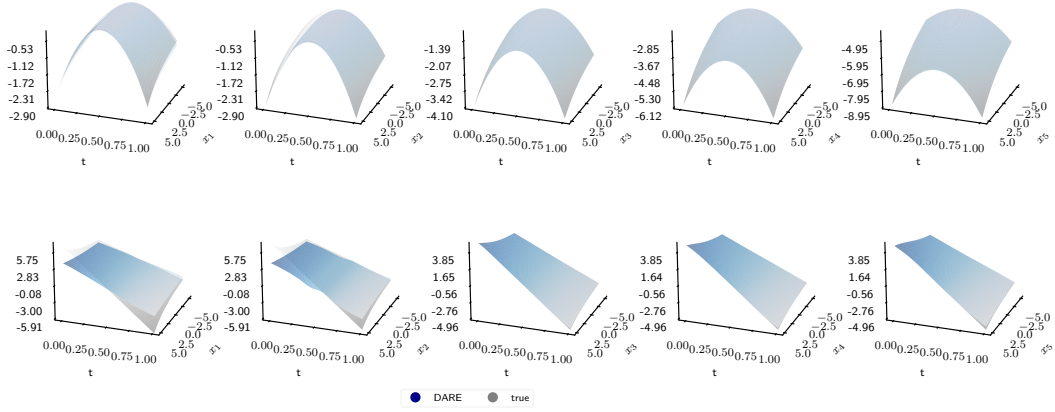


Figure 10: True and approximated (with DARE) value function (31) for $t \in [0, T]$ and $\mathbf{X} = X_1, X_2, X_3, X_4, X_5 \in [-5, 5]^5$. Each surface corresponds to the value function for time and one dimension in X , where the value of the system in all other dimensions is fixed to $x_i = 0$.

Gaussian measurement noise $\epsilon_i \sim \mathcal{N}(0, s^2)$ for $s > 0$. GPs are fully specified by a mean function $\mu : \mathcal{X} \mapsto \mathbb{R}$ and a covariance (kernel) function $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$. In particular, if $f \sim \mathcal{GP}(\mu, k)$ and \mathbf{X}_* is a set of test points in the domain \mathcal{X} of the GP, then the set of random variables \mathbf{f}_* is Gaussian with parameters $\mathcal{N}(\boldsymbol{\mu}_*, \mathbf{K}_{*,*})$, where

$$\boldsymbol{\mu}_* = \{\mu(\mathbf{x})\}_{\mathbf{x} \in \mathbf{X}_*} \quad \text{and} \quad \mathbf{K}_{*,*} = \{k(\mathbf{x}, \mathbf{x}')\}_{(\mathbf{x}, \mathbf{x}') \in \mathbf{X}_*}.$$

681 A convenient property of GPs is that one computes the posterior distribution with analytic formulae.
 682 Suppose we collect n noisy observations $\mathbf{y} = \{y_1, \dots, y_n\}$ at the domain points $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$,
 683 where $y_i = f(\mathbf{x}_i) + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, s^2)$. Then, the posterior distribution over f given the previous
 684 (training) observations \mathbf{X} and \mathbf{y} , is also a GP with mean function μ_{post} and covariance function k_{post}
 685 given by

$$\begin{cases} \mu_{\text{post}}(\mathbf{x}_*) &= \mathbf{k}(\mathbf{x}_*, \mathbf{X}) (\mathbf{K} + s^2 \mathbf{I})^{-1} \mathbf{y}, \\ k_{\text{post}}(\mathbf{x}_*, \mathbf{x}'_*) &= k(\mathbf{x}_*, \mathbf{x}'_*) \\ &\quad - \mathbf{k}(\mathbf{x}_*, \mathbf{X}) (\mathbf{K} + s^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}'_*), \end{cases} \quad (41)$$

where

$$\mathbf{k}(\mathbf{x}_*, \mathbf{X}) = \mathbf{k}(\mathbf{X}, \mathbf{x}_*)^\top = (k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_n))$$

686 is the n -dimensional covariance vector of the test point \mathbf{x}_* with training points $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$,
 687 $\mathbf{K} = (k(\mathbf{x}_i, \mathbf{x}_j))_{i,j \in \{1, \dots, n\}}$ is the positive semi-definite kernel matrix from training data, and \mathbf{I} is
 688 the n -dimensional identity matrix. See Figure I for an example.

689 Let the elements of the vector $\boldsymbol{\theta} \in \Theta$ be hyper-parameters of the prior's kernel function and s^2 is the
 690 variance of the i.i.d. Gaussian noise that corrupts reward observations. Both $\boldsymbol{\theta}$ and s^2 are inferred
 691 with the log marginal likelihood of the data given by

$$\begin{aligned} L(\boldsymbol{\theta}, s) &= \log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}, s) & (42) \\ &= -\frac{1}{2} \log \left(\det(\mathbf{K}_\boldsymbol{\theta} + s^2 \mathbf{I}) \right) \\ &\quad - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_\boldsymbol{\theta} + s^2 \mathbf{I})^{-1} \mathbf{y} - \frac{n}{2} \log(2\pi), \end{aligned}$$

692 for a zero-mean GP, where \mathbf{X} and \mathbf{y} are the n training samples and $\mathbf{K}_\boldsymbol{\theta}$ is the prior's positive
 693 covariance matrix with kernel $k_\boldsymbol{\theta}$. The vector of hyper-parameters $\boldsymbol{\theta}$ and the variance s^2 maximize
 694 the quantity (42), i.e., $(\boldsymbol{\theta}^*, s^*) \in \arg \max_{\boldsymbol{\theta} \in \Theta, s \in \mathbb{R}^+} L(\boldsymbol{\theta}, s)$, which one solves with classical gradient descent-

695 based optimization algorithms.

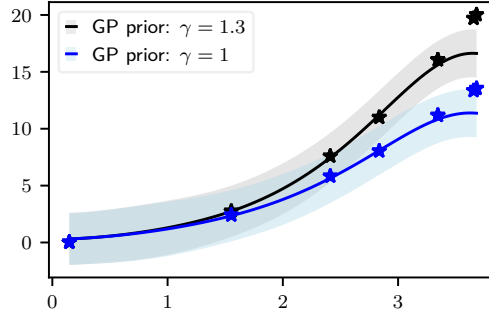


Figure 11: Two GPs fitted to $f(u) = u^{1+\gamma_i}$ for $\gamma_0 = 1.3$ and $\gamma_1 = 1$.

696 J Reinforcement Learning Benchmarks

697 Table 1 shows the average elapsed real time i.e. wall-clock time comparison of Reinforcement
 698 Learning Algorithms (PPO, SAC, A2C) with that of DARE rounded to the nearest second training
 699 on the LQG problem from section 4.3. We utilize the implementations in the Stable-Baselines-3
 700 library [41] for the Reinforcement Learning experiments.

701 We measure the elapsed time from the first step in the training loop until convergence using the `time`
 702 module from Python 3.9.19.

703 The hyperparameter ranges in Table 5 are determined in a way that includes the recommended values
 704 on the Stable-Baselines-3 [41] implementation used in this work.

Table 5: Hyperparameter Range considered for each RL algorithm

Hyperparameter	Algorithm	Range
Learning Rate	PPO	$[10^{-4}, 10^{-3}]$ (continuous range)
	SAC	$[10^{-4}, 10^{-3}]$ (continuous range)
	A2C	$[10^{-4}, 10^{-3}]$ (continuous range)
Action Repeats	PPO	{10, 20, 40, 50, 100}
	SAC	{10, 20, 40, 50, 100}
	A2C	{10, 20, 40, 50, 100}
Number of Steps	PPO	[128, 4096] (integer range)
	A2C	[4, 32] (integer range)

705 K Impact Statement

706 This paper presents a novel deep learning methodology for solving decision-making problems in
 707 noisy and non-stationary environments, with wide-ranging applications in finance, robotics, and
 708 biology. Our contribution is a highly accurate and efficient method for solving model predictive
 709 control problems. Possible implications include more efficient and effective risk management in
 710 finance, safer robot-human interaction, and improved biomedical engineering. We use tractable
 711 examples to test our approach and to demonstrate that our model produces reasonable policies.
 712 Before implementing our model for critical problems, we believe further specific experimentation
 713 and validation is necessary.

714 **NeurIPS Paper Checklist**

715 **1. Claims**

716 Question: Do the main claims made in the abstract and introduction accurately reflect the
717 paper's contributions and scope?

718 Answer: [\[Yes\]](#)

719 Justification: The claims in the introduction and abstract accurately reflect the set of problems
720 the paper is focused on and how the proposed methods were tested and compared.

721 Guidelines:

- 722 • The answer NA means that the abstract and introduction do not include the claims
723 made in the paper.
- 724 • The abstract and/or introduction should clearly state the claims made, including the
725 contributions made in the paper and important assumptions and limitations. A No or
726 NA answer to this question will not be perceived well by the reviewers.
- 727 • The claims made should match theoretical and experimental results, and reflect how
728 much the results can be expected to generalize to other settings.
- 729 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
730 are not attained by the paper.

731 **2. Limitations**

732 Question: Does the paper discuss the limitations of the work performed by the authors?

733 Answer: [\[Yes\]](#)

734 Justification: Yes, we discuss the limitations of the work like its requirements for PDE
735 representations that are not required by other fields like Reinforcement Learning that also
736 deal with high dimensional control problems.

737 Guidelines:

- 738 • The answer NA means that the paper has no limitation while the answer No means that
739 the paper has limitations, but those are not discussed in the paper.
- 740 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 741 • The paper should point out any strong assumptions and how robust the results are to
742 violations of these assumptions (e.g., independence assumptions, noiseless settings,
743 model well-specification, asymptotic approximations only holding locally). The authors
744 should reflect on how these assumptions might be violated in practice and what the
745 implications would be.
- 746 • The authors should reflect on the scope of the claims made, e.g., if the approach was
747 only tested on a few datasets or with a few runs. In general, empirical results often
748 depend on implicit assumptions, which should be articulated.
- 749 • The authors should reflect on the factors that influence the performance of the approach.
750 For example, a facial recognition algorithm may perform poorly when image resolution
751 is low or images are taken in low lighting. Or a speech-to-text system might not be
752 used reliably to provide closed captions for online lectures because it fails to handle
753 technical jargon.
- 754 • The authors should discuss the computational efficiency of the proposed algorithms
755 and how they scale with dataset size.
- 756 • If applicable, the authors should discuss possible limitations of their approach to
757 address problems of privacy and fairness.
- 758 • While the authors might fear that complete honesty about limitations might be used by
759 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
760 limitations that aren't acknowledged in the paper. The authors should use their best
761 judgment and recognize that individual actions in favor of transparency play an impor-
762 tant role in developing norms that preserve the integrity of the community. Reviewers
763 will be specifically instructed to not penalize honesty concerning limitations.

764 **3. Theory Assumptions and Proofs**

765 Question: For each theoretical result, does the paper provide the full set of assumptions and
766 a complete (and correct) proof?

767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820

Answer: [Yes]

Justification: Section E of the Appendix includes the proofs with the relevant definitions and propositions.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide all experimental details along with an anonymized repository containing all code.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

821 Question: Does the paper provide open access to the data and code, with sufficient instruc-
822 tions to faithfully reproduce the main experimental results, as described in supplemental
823 material?

824 Answer: [Yes]

825 Justification: All code can be run using open source packages without external data.

826 Guidelines:

- 827 • The answer NA means that paper does not include experiments requiring code.
- 828 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
829 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 830 • While we encourage the release of code and data, we understand that this might not be
831 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
832 including code, unless this is central to the contribution (e.g., for a new open-source
833 benchmark).
- 834 • The instructions should contain the exact command and environment needed to run to
835 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
836 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 837 • The authors should provide instructions on data access and preparation, including how
838 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 839 • The authors should provide scripts to reproduce all experimental results for the new
840 proposed method and baselines. If only a subset of experiments are reproducible, they
841 should state which ones are omitted from the script and why.
- 842 • At submission time, to preserve anonymity, the authors should release anonymized
843 versions (if applicable).
- 844 • Providing as much information as possible in supplemental material (appended to the
845 paper) is recommended, but including URLs to data and code is permitted.

846 6. Experimental Setting/Details

847 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
848 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
849 results?

850 Answer: [Yes]

851 Justification: Outlines of each of the experiments as well as our choice of hyperparameters
852 is given in Section 4.1 and Appendix J.

853 Guidelines:

- 854 • The answer NA means that the paper does not include experiments.
- 855 • The experimental setting should be presented in the core of the paper to a level of detail
856 that is necessary to appreciate the results and make sense of them.
- 857 • The full details can be provided either with the code, in appendix, or as supplemental
858 material.

859 7. Experiment Statistical Significance

860 Question: Does the paper report error bars suitably and correctly defined or other appropriate
861 information about the statistical significance of the experiments?

862 Answer: [Yes]

863 Justification: Each test is run with a sufficient across a sufficient number of seeds when
864 necessary.

865 Guidelines:

- 866 • The answer NA means that the paper does not include experiments.
- 867 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
868 dence intervals, or statistical significance tests, at least for the experiments that support
869 the main claims of the paper.
- 870 • The factors of variability that the error bars are capturing should be clearly stated (for
871 example, train/test split, initialization, random drawing of some parameter, or overall
872 run with given experimental conditions).

- 873 • The method for calculating the error bars should be explained (closed form formula,
874 call to a library function, bootstrap, etc.)
- 875 • The assumptions made should be given (e.g., Normally distributed errors).
- 876 • It should be clear whether the error bar is the standard deviation or the standard error
877 of the mean.
- 878 • It is OK to report 1-sigma error bars, but one should state it. The authors should
879 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
880 of Normality of errors is not verified.
- 881 • For asymmetric distributions, the authors should be careful not to show in tables or
882 figures symmetric error bars that would yield results that are out of range (e.g. negative
883 error rates).
- 884 • If error bars are reported in tables or plots, The authors should explain in the text how
885 they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

887 Question: For each experiment, does the paper provide sufficient information on the com-
888 puter resources (type of compute workers, memory, time of execution) needed to reproduce
889 the experiments?

890 Answer: [Yes]

891 Justification: We provide hardware information in Appendix F. Our tests do not require any
892 specialized or intensive computing power.

893 Guidelines:

- 894 • The answer NA means that the paper does not include experiments.
- 895 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
896 or cloud provider, including relevant memory and storage.
- 897 • The paper should provide the amount of compute required for each of the individual
898 experimental runs as well as estimate the total compute.
- 899 • The paper should disclose whether the full research project required more compute
900 than the experiments reported in the paper (e.g., preliminary or failed experiments that
901 didn't make it into the paper).

9. Code Of Ethics

903 Question: Does the research conducted in the paper conform, in every respect, with the
904 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

905 Answer: [Yes]

906 Justification: We conform with the NeurIPS Code of Ethics.

907 Guidelines:

- 908 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 909 • If the authors answer No, they should explain the special circumstances that require a
910 deviation from the Code of Ethics.
- 911 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
912 eration due to laws or regulations in their jurisdiction).

10. Broader Impacts

914 Question: Does the paper discuss both potential positive societal impacts and negative
915 societal impacts of the work performed?

916 Answer: [Yes]

917 Justification: An impact statement is included in Appendix K.

918 Guidelines:

- 919 • The answer NA means that there is no societal impact of the work performed.
- 920 • If the authors answer NA or No, they should explain why their work has no societal
921 impact or why the paper does not address societal impact.

- 922
- 923
- 924
- 925
- 926
- 927
- 928
- 929
- 930
- 931
- 932
- 933
- 934
- 935
- 936
- 937
- 938
- 939
- 940
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

941 11. Safeguards

942 Question: Does the paper describe safeguards that have been put in place for responsible
943 release of data or models that have a high risk for misuse (e.g., pretrained language models,
944 image generators, or scraped datasets)?

945 Answer: [NA]

946 Justification: We believe our model does not have a high risk for misuse.

947 Guidelines:

- 948
- 949
- 950
- 951
- 952
- 953
- 954
- 955
- 956
- 957
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

958 12. Licenses for existing assets

959 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
960 the paper, properly credited and are the license and terms of use explicitly mentioned and
961 properly respected?

962 Answer: [Yes]

963 Justification: All sources are cited and we use publicly available packages.

964 Guidelines:

- 965
- 966
- 967
- 968
- 969
- 970
- 971
- 972
- 973
- 974
- 975
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 976
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 977
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.
- 978
- 979

980 **13. New Assets**

981 Question: Are new assets introduced in the paper well documented and is the documentation
982 provided alongside the assets?

983 Answer: [NA]

984 Justification: The paper does not release new assets.

985 Guidelines:

- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 986
- 987
- 988
- 989
- 990
- 991
- 992
- 993

994 **14. Crowdsourcing and Research with Human Subjects**

995 Question: For crowdsourcing experiments and research with human subjects, does the paper
996 include the full text of instructions given to participants and screenshots, if applicable, as
997 well as details about compensation (if any)?

998 Answer: [NA]

999 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1000 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 1001
- 1002
- 1003
- 1004
- 1005
- 1006
- 1007
- 1008

1009 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
1010 Subjects**

1011 Question: Does the paper describe potential risks incurred by study participants, whether
1012 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1013 approvals (or an equivalent approval/review based on the requirements of your country or
1014 institution) were obtained?

1015 Answer: [NA]

1016 Justification: The paper does not involve crowdsourcing nor research with human subjects

1017 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
 - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
 - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.
- 1018
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025
- 1026
- 1027