# Approximate Optimal Online Continuous-Time Path-Planner with Static Obstacle Avoidance

Patrick Walters, Rushikesh Kamalapurkar, and Warren E. Dixon

*Abstract*—Online approximation of the optimal path for a control affine nonlinear autonomous agent subject to input and state constraints (e.g., actuator saturation, obstacles, no-enter zones) is considered. A model-based adaptive dynamic programming technique is implemented to locally estimate the unknown value function associated with the optimal path-planning problem. By performing a local approximation, the locations of the static obstacles do not need to be known until the obstacles are within a defined approximation window. The developed feedback policy guarantees ultimately bounded convergence of the approximated path to the optimal path without the requirement of persistence of excitation, typically required for online adaptive dynamic programming. Simulation results are presented to illustrate the performance of the proposed method.

## I. INTRODUCTION

Path-planning approaches can be divided into two types, pregenerative and reactive [1]. Pregenerative methods compute a path before a mission begins (c.f [1]–[4]), while reactive methods determine a path as the agent progresses through its environment. Planning optimal paths is of particular interest. Often pregenerative methods are used in planning optimal paths, which in the event of a disturbance requires the agent to take a non-optimal trajectory to return to the original optimal path or continually execute the planner. In contrast, an optimal reactive method (feedback motion planner) has the advantage of generating a policy that provides optimal feedback if the agent is forced off its original path.

In developing an optimal path for autonomous agents, it is often necessary to consider the agent's dynamics. In general, it is difficult to develop optimal path-planing strategies for nonlinear dynamics. One method of dealing with the challenges of path-planning under differential constraints is to pose the problem as an of optimal control problem. The corresponding Hamilton-Jacobi-Bellmen (HJB) equation can be numerically approximated to produce feedback policies. Dynamic programming with interpolation has been used as a feedback motion planner to compute an approximate optimal path through value iteration in results such as [5]. However, similar to pregenerative graph search methods (e.g., A*, Dijkstra), difficulties arise related to the state discretization

as the order of the dynamics increase [6]. In results such as [7] and [8], feedback-based path-planning is generated offline by solving the HJB numerically. In the event of a change in the environment, such results would be required to recalculate a new optimal plan offline. In this result, we consider adaptive dynamic programming (ADP), which has been used to approximate the solution to HJB equation of general nonlinear systems online using parametric function approximation techniques (e.g., [9]–[13]).

Further complicating the task of optimal path-planning, state constraints (e.g., obstacles, no-enter zones) are often present en route to an objective. Autonomous agents are not able to sense all obstacles a priori, e.g., obstacles may remain undiscovered until they fall within a given sensing range. A recent advance in ADP bases the parametric approximation of the solution to the HJB on state-following (StaF) kernels [14]. These StaF kernels yield a local approximation of the HJB around the current state. By only utilizing information near the current state to approximate the solution, StaF does not require knowledge of obstacles outside an approximation window.

In addition to state constraints, input constraints inherent to the agent (e.g., maximum speed) are also important to consider. Results such as [12], [13], [15] have considered input constraints within the ADP framework. Utilizing a generalized non-quadratic local cost [16], the results in [12], [13], [15] yield a bounded approximate optimal controller.

Inspired by the advances in [12]–[15], an optimal feedback-based path-planner is developed in this paper that respects input and state constraints. The developed planner differs from optimal controllers in that it tackles the challenges specific to obstacle avoidance within the path-planning strategy. The local approximation in [14] enables the handling of obstacles not known a priori, but also introduces time-varying parameters in the parametric representation of the solution to the HJB. The time-varying parameters cause the estimation of the parameters to be in a near constant transient state making it difficult to prove that the generated feedback policy avoids the obstacles. This technical challenge motivates the introduction of an auxiliary feedback term to assist in navigating an agent around obstacles, and a scheduling function to switch between the approximate optimal feedback plan and the auxiliary feedback plan. Switching to the auxiliary feedback plan when the agent risks hitting an obstacle ensures obstacle avoidance. The proposed model-based ADP method approximates the optimal path using a combination of on-policy and off-policy

data, eliminating the need for physical exploration of the state space often associated with ADP. A Lyapunov-based stability analysis is presented which guarantees ultimately bounded convergence of the approximated path to the optimal path. Simulation results validate the proposed method.

## II. ONLINE APPROXIMATE OPTIMAL PATH-PLANNING

The objective of the proposed path-planner is to optimally navigate to a setpoint while avoiding static obstacles and adhering to the agent's actuation constraints.

### A. Problem Formulation

Consider a nonlinear control affine dynamical system of the form

$$\dot{\zeta}(t) = \overline{f}(\zeta(t)) + \overline{g}(\zeta(t))\,\overline{u}(t), \tag{1}$$

$t \in [t_0, \infty)$, where $t_0$ denotes the initial time, $\zeta : [t_0, \infty) \to \mathbb{R}^n$ denotes the system state, $\overline{f} : \mathbb{R}^n \to \mathbb{R}^n$ denotes the drift dynamics, $\overline{g} : \mathbb{R}^n \to \mathbb{R}^{n \times m}$ denotes the control effectiveness, and $\overline{u} : [t_0, \infty) \to \mathbb{R}^m$ denotes the control input.

The path-planning problem may be posed as a constrained infinite-horizon nonlinear regulation problem, i.e., to design a control signal $\overline{u}$ to minimize a subsequently defined cost function subject to the dynamic constraint in (1), while avoiding the obstacles and obeying $\sup_t (\overline{u}_i) \le u_{sat} \ \forall i = 1, \dots, m$, where $\overline{u} = [\overline{u}_1, \cdots, \overline{u}_m]^T$ and $u_{sat} \in [0, \infty)$ is the control saturation constant.

Static obstacles represent hard constraints that must be taken into account in the development of the approximate optimal path-planner. To this end, an auxiliary controller is subsequently developed to assist in navigating an agent around obstacles. The auxiliary controller is denoted by $u_s : \mathbb{R}^n \to \mathbb{R}^m$, where $u_s = [u_{s_1}, \cdots, u_{s_m}]^T$ and $\sup_t (u_{s_i}) \le u_{sat} \ \forall i = 1, 2, \dots, m$. To facilitate the development of $u_s$, obstacles are augmented with a perimeter that extends from their borders denoting an unsafe region as illustrated in Figure 1. A smooth scheduling function $s : \mathbb{R}^n \to [0, a]$, where $a < 1$, is used to transition between the approximate optimal controller $u : [t_0, \infty) \to \mathbb{R}^m$ and the auxiliary controller $u_s$ without introducing discontinuities to the system dynamics. The scheduling function and auxiliary controller are designed such that they are functions of the state and that they drive all state trajectories away from obstacles. In Section IV, the auxiliary controller $u_s$ and the scheduling function $s$ are designed for a specific system.

The control input $\overline{u}$ is defined as

$$\overline{u}(\zeta, t) \triangleq s(\zeta)\,u_s(\zeta) + (1 - s(\zeta))\,u(t), \tag{2}$$

where $u$ is the subsequently designed approximate optimal controller. Based on (2), the agent dynamics can be rewritten as

$$\dot{\zeta}(t) = f(\zeta(t)) + g(\zeta(t))\,u(t), \tag{3}$$

where $f : \mathbb{R}^n \to \mathbb{R}^n$ denotes the augmented drift dynamics defined as

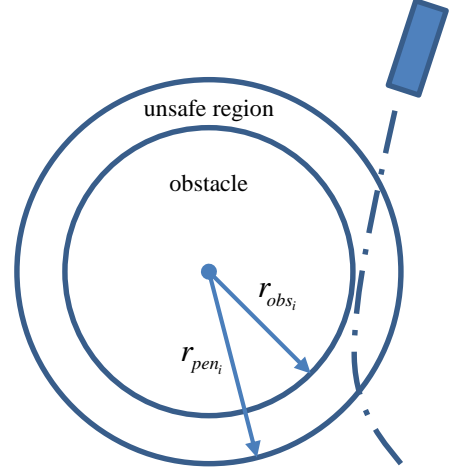$$f(\zeta) \triangleq \overline{f}(\zeta) + \overline{g}(\zeta)\,s(\zeta)\,u_s(\zeta)$$



Figure 1. Obstacles are augmented with an unsafe region that extends from it's border.

and $g : \mathbb{R}^n \to \mathbb{R}^{n \times m}$ denotes the augmented control effectiveness defined as

$$g(\zeta) \triangleq \overline{g}(\zeta)\,(1 - s(\zeta)).$$

To account for actuator saturation and the unsafe regions, the cost function is defined as

$$J(\zeta, u) \triangleq \int_{t_o}^{\infty} r(\zeta(\tau), u(\tau))\, d\tau, \tag{4}$$

where $r : \mathbb{R}^n \times \mathbb{R}^m \to [0, \infty)$ is the local cost defined as

$$r(\zeta, u) \triangleq \zeta^T Q \zeta + P(\zeta) + U(u), \tag{5}$$

subject to the dynamic constraint in (3), where $Q \in \mathbb{R}^{n \times n}$ is a constant, user defined, symmetric positive definite weighting matrix, $P : \mathbb{R}^n \to \mathbb{R}$ is a non-negative continuous function penalizing trajectories that enter the unsafe region, and $U : \mathbb{R}^m \to \mathbb{R}$ is a positive definite function penalizing control effort. The matrix $Q$ has the property $\underline{q}\,\|\xi_q\|^2 \le \xi_q^T Q \xi_q \le \overline{q}\,\|\xi_q\|^2$, $\forall \xi_q \in \mathbb{R}^n$ where $\underline{q}$ and $\overline{q}$ are positive constants.

The positive definite function $U$ in (5) is defined as [12], [17]

$$U(u) \triangleq 2 \sum_{i=1}^{m} \left[ \int_0^{u_i} \left( u_{sat} r_i \tanh^{-1}\left( \frac{\xi_{u_i}}{u_{sat}} \right) \right) d\xi_{u_i} \right] \tag{6}$$

where $u_i$ is the $i^{th}$ element of $u$, $\xi_{u_i}$ is an integration variable, and $R \in \mathbb{R}^{m \times m}$ is a diagonal positive definite, user defined, weighting matrix given as $R = \text{diag}\,([r_1, r_2, \cdots, r_m])$.

The infinite-time scalar value function $V : \mathbb{R}^n \to [0, \infty)$ for the optimal solution is written as

$$V(\zeta) = \min_u \int_{t_0}^{\infty} r(\zeta(\tau), u(\tau))\, d\tau. \tag{7}$$

The objective of the optimal path-planner is to find the optimal policy $u^* : \mathbb{R}^n \to \mathbb{R}^m$ that minimizes the performance index (4) with the local cost (5) subject to the dynamic constraint in (3).

The optimal value function is characterized by the HJB, which is given as

$$\frac{\partial V(\zeta)}{\partial \zeta}\left(f(\zeta) + g(\zeta)u^*(\zeta)\right) + r(\zeta, u^*(\zeta)) = 0 \quad (8)$$

with the boundary condition $V(0) = 0$. The optimal control policy can be determined from (8) as

$$u^*(\zeta) = -u_{sat} \tanh\left(\frac{1}{2u_{sat}}R^{-1}g^T\left(\frac{\partial V(\zeta)}{\partial \zeta}\right)^T\right). \quad (9)$$

The analytical expression for the optimal path in (9) requires knowledge of the value function which is the solution to the HJB equation in (8). The HJB equation is a partial differential equation which is generally infeasible to solve analytically; hence, an approximate solution is sought.

### B. Local Approximation of the Value Function

The subsequent development is based on an approximation of the value function and optimal policy. Differing from previous ADP literature (e.g., [?], [10], [12], [13], [18]) that seeks a global policy, the following development seeks only a local policy. Instead of generating an approximation of the value function over the entire operating region, we aim to approximate a small region about the current state. With the region of approximation limited to a small range about the current state, one only needs to assume that there may exist an obstacle or obstacles outside the local approximation. Once inside the local approximation window, the optimal policy will adapt to avoid the obstacle. Despite the uncertainty of distant obstacles, the following development yields guaranteed stability of the state and convergence to the optimal path.

Leveraging the results of [14], StaF kernels are employed to approximate the local policy on some small compact set $B_r(\zeta)$, i.e., the approximation window, around the state $\zeta$. The StaF representation of the value function and optimal policy are given as

$$V(\zeta) = W(\zeta)^T \sigma(\zeta, c(\zeta)) + \epsilon(\zeta),$$

$$u^*(\zeta) = -u_{sat} \tanh\left(\frac{R^{-1}g(\zeta)^T}{2u_{sat}}\left(\sigma'(\zeta, c(\zeta))^T W(\zeta)\right.\right.$$
$$\left.\left. + \epsilon'(\zeta)\right)\right), \quad (10)$$

respectively, where $W : \mathbb{R}^n \to \mathbb{R}^l$ is the ideal weight vector, $\sigma : \mathbb{R}^n \to \mathbb{R}^l$ is a continuously differentiable kernel function, and $\epsilon : \mathbb{R}^n \to \mathbb{R}$ is the continuously differentiable function reconstruction error, and $\sigma' : \mathbb{R}^n \to \mathbb{R}^{l \times n}$ and $\epsilon' : \mathbb{R}^n \to \mathbb{R}^n$ are the partial derivatives with respect to the state. Note that the centers of the kernel function change as the system state changes; therefore, the ideal weight vector $W$ is a time-varying function. The approximations of the value function and the optimal policy are defined as

$$\hat{V}\left(\zeta, \hat{W}_c\right) \triangleq \hat{W}_c^T \sigma(\zeta, c(\zeta)), \quad (11)$$

$$\hat{u}\left(\zeta, \hat{W}_a\right) \triangleq -u_{sat} \tanh\left(\frac{R^{-1}g(\zeta)^T}{2u_{sat}}\sigma'(\zeta, c(\zeta))^T \hat{W}_a\right), \quad (12)$$

where $c(\zeta) \in B_r(\zeta)$ is the StaF kernel center, and $\hat{W}_c, \hat{W}_a \in \mathbb{R}^l$ are estimates of the ideal weight vector $W$.

Substituting the approximations from (11) and (12) into (8), results in a residual error $\delta : \mathbb{R}^n \times \mathbb{R}^l \times \mathbb{R}^l \to \mathbb{R}$ called the Bellman error given by

$$\delta\left(\zeta, \hat{W}_c, \hat{W}_a\right) = r\left(\zeta, \hat{u}\left(\zeta, \hat{W}_a\right)\right)$$
$$+ \frac{\partial \hat{V}\left(\zeta, \hat{W}_c\right)}{\partial \zeta}\left(f(\zeta) + g(\zeta)\hat{u}\left(\zeta, \hat{W}_a\right)\right). \quad (13)$$

### C. Online Learning

At a given time instant $t$, the Bellman error $\delta_t : [0, \infty) \to \mathbb{R}$ is

$$\delta_t(t) = \delta\left(\zeta(t), \hat{W}_c(t), \hat{W}_a(t)\right),$$

where $\hat{W}_c(t)$ and $\hat{W}_a(t)$ denote the estimates of the ideal weights at time $t$, and $\zeta(t)$ denotes the state of the system in (1) at time $t$ starting from initial time $t_o$ and initial state $\zeta_o$ under the influence of the state feedback controller

$$u(t) = \hat{u}\left(\zeta(t), \hat{W}_a(t)\right).$$

The Bellman error is extrapolated to off-policy sampled states $\{\zeta_k(t) \in B_r(\zeta(t)) | k = 1, 2, \ldots, N\}$ that follow the system state. The extrapolated Bellman error $\delta_k : [0, \infty) \to \mathbb{R}$ is given as

$$\delta_k(t) = \hat{W}_c^T(t)w_k(t) + r(\zeta_k(t), \hat{u}_k(t)),$$

where

$$\hat{u}_k(t) = -u_{sat}\tanh\left(\frac{R^{-1}g}{2u_{sat}}(\zeta_k(t))^T\right.$$
$$\left.\sigma'(\zeta_k(t), c(\zeta(t)))^T \hat{W}_a(t)\right)$$

and

$$\omega_k(t) = \sigma'(\zeta_k(t), c(\zeta(t)))\left[f(\zeta_k(t))\right.$$
$$\left. + g(\zeta_k(t))\hat{u}\left(\zeta(t), \hat{W}_a(t)\right)\right].$$

The value function least squares update law based on the minimization of the instantaneous and extrapolated Bellman error is given by

$$\dot{\hat{W}}_c(t) = -\Gamma(t)\left(k_{c1}\frac{\partial \delta_t(t)}{\partial \hat{W}_c}\frac{\delta_t(t)}{\rho(t)}\right.$$
$$\left. + \frac{k_{c2}}{N}\sum_{k=1}^{N}\frac{\partial \delta_k(t)}{\partial \hat{W}_c}\frac{\delta_k(t)}{\rho_k(t)}\right), \quad (14)$$

$$\dot{\Gamma}(t) = \begin{cases} \beta\Gamma(t) - k_{c1}\Gamma(t)\frac{\omega(t)\omega(t)^T}{\rho(t)}\Gamma(t), & \|\Gamma(t)\| \leq \overline{\Gamma} \\ 0, & \text{otherwise} \end{cases},$$

(15)

where $k_{c1}$, $k_{c2} \in \mathbb{R}$ are positive adaptation gains, $\|\Gamma(t_0)\| = \|\Gamma_0\| \leq \overline{\Gamma}$ is the initial adaptation gain, $\overline{\Gamma} \in \mathbb{R}$ is a positive saturation gain, $\beta \in \mathbb{R}$ is a positive forgetting factor,

$$\omega(t) = \sigma'(\zeta(t), c(\zeta(t)))[f(\zeta(t)) \\ + g(\zeta(t))\,\hat{u}(\zeta(t), \hat{W}_a(t))]$$

is the instantaneous regressor matrix,

$$\rho(t) = 1 + k_\rho\omega(t)^T\omega(t)$$

is the instantaneous normalization constant,

$$\rho_k(t) = 1 + k_\rho\omega_k(t)^T\omega_k(t)$$

is the extrapolated normalization constant, and $k_\rho \in \mathbb{R}$ is a positive gain.

The policy update law is given by

$$\dot{\hat{W}}_a(t) = \text{proj}\left\{-k_a\left(\hat{W}_a(t) - \hat{W}_c(t)\right)\right\},$$

(16)

where $k_a \in \mathbb{R}$ is an positive gain, and $\text{proj}\{\cdot\}$ is a smooth projection operator[1] used to bound the weight estimates. The weight estimation errors are then defined as $\tilde{W}_c(t) \triangleq W(\zeta(t)) - \hat{W}_c(t)$ and $\tilde{W}_a(t) \triangleq W(\zeta(t)) - \hat{W}_a(t)$.

In Section III, Bellman error extrapolation is employed to establish ultimately bounded convergence of the approximate policy to the optimal policy without requiring persistence of excitation provided the following assumption is satisfied.

**Assumption 1.** There exists a strictly positive constant $\underline{c}$ such that

$$\underline{c} = \inf_{t \in [t_0, \infty)}\left[\lambda_{min}\left(\sum_{k=1}^{N}\frac{\omega_k(t)\omega_k(t)^T}{\rho_k(t)^2}\right)\right].$$

In general, Assumption 1 cannot be guaranteed to hold a priori; however, heuristically, the condition can be met by sampling redundant data, i.e., $N \gg l$.

## III. Stability Analysis

For notational brevity, all function dependencies from previous sections are henceforth suppressed. Let the notation $F_k$ denote the function $F(\zeta, \cdot)$ evaluated at the sampled state, i.e., $F_k(\cdot) = F(\zeta_k, \cdot)$. An unmeasurable form of the Bellman error can be written as

$$\delta = -\tilde{W}_c^T\omega + W^T\sigma'g(\hat{u} - u^*) + U(\hat{u}) - U(u^*) \\ - \epsilon'(f + gu^*).$$

(17)

Similarly, the Bellman error at the extrapolated points can be written as

$$\delta_k = -\tilde{W}_c^T\omega_k + W^T\sigma_k'g_k(\hat{u}_k - u_k^*) + U(\hat{u}_k) - U(u_k^*) \\ - \epsilon_k'(f_k + g_ku_k^*).$$

(18)

[1] See Section 4.4 in [19] or Remark 3.6 in [20] for details of the projection operator.

To facilitate the subsequent stability analysis, consider the candidate Lyapunov function $V_L : \mathbb{R}^n \times \mathbb{R}^l \times \mathbb{R}^l \to [0, \infty)$ given as

$$V_L(Z) = V + \frac{1}{2}\tilde{W}_c^T\Gamma^{-1}\tilde{W}_c + \frac{1}{2}\tilde{W}_a^T\tilde{W}_a,$$

where $Z \triangleq \begin{bmatrix} \zeta^T & \tilde{W}_c^T & \tilde{W}_a^T \end{bmatrix}^T \in \chi \times \mathbb{R}^l \times \mathbb{R}^l$. Since the value function $V$ in (7) is positive definite [21] using Lemma 4.3 of [22], $V_L$ can be bounded by

$$\underline{v_L}(\|Z\|) \leq V_L(Z) \leq \overline{v_L}(\|Z\|),$$

(19)

where $\underline{v_L}$, $\overline{v_L} : [0, \infty) \to [0, \infty)$ are class $\mathcal{K}$ functions. Define the constant $K \in \mathbb{R}$ as

$$K \triangleq \sqrt{\frac{\iota_c^2}{\alpha(2\underline{c} - k_a)} + \frac{\iota_a^2}{\alpha k_a} + \frac{\iota}{\alpha}}$$

where $\alpha \triangleq \min\left\{\frac{q}{2}, \left(\frac{\underline{c}}{2} - \frac{k_a}{4}\right), \frac{k_a}{4}\right\}$ and $\iota_c$, $\iota_a$, $\iota$ are positive constants.

**Theorem 1.** *Provided Assumption 1 is satisfied along with the sufficient conditions*

$$\underline{c} > \frac{k_a}{2},$$

(20)

$$K < \underline{v_L}^{-1}(\overline{v_L}(r)),$$

(21)

*where $r \in \mathbb{R}$ is the radius of the compact set $\beta_L$, then the policy in (12) with the update laws in (14)-(16) guarantee ultimately bounded regulation of the state $\zeta$ and of the approximated policies $\hat{u}$ to the optimal policy $u^*$.*

*Proof:* The time derivative of the candidate Lyapunov function is

$$\dot{V}_L = \frac{\partial V}{\partial \zeta}f + \frac{\partial V}{\partial \zeta}g\hat{u} - \frac{1}{2}\tilde{W}_c^T\Gamma^{-1}\dot{\Gamma}\Gamma^{-1}\tilde{W}_c \\ - \tilde{W}_c^T\Gamma^{-1}\left(\dot{W} - \dot{\hat{W}}_c\right) - \tilde{W}_a^T\left(\dot{W} - \dot{\hat{W}}_a\right).$$

(22)

Using Theorem 2 in [14], the time derivative of the ideal weights can be expressed as

$$\dot{W} = W'(f + g\hat{u}).$$

(23)

Substituting (14)-(18) and (23) yields

$$\dot{V}_L = \frac{\partial V}{\partial \zeta}f + \frac{\partial V}{\partial \zeta}g\hat{u} \\ - \frac{1}{2}\tilde{W}_c^T\Gamma^{-1}\left[\left(\left(\beta\Gamma - k_{c1}\Gamma\frac{\omega\omega^T}{\rho}\Gamma\right)\mathbf{1}_{\|\Gamma\| \leq \overline{\Gamma}}\right]\Gamma^{-1}\tilde{W}_c \\ + \tilde{W}_c^T\left[k_{c1}\frac{\omega}{\rho}\delta + \frac{k_{c2}}{N}\sum_{j=1}^{N}\frac{\omega_k}{\rho_k}\delta_k\right] - \tilde{W}_c^T\Gamma^{-1}W'(f + g\hat{u}) \\ + \tilde{W}_a^Tk_a\left(\hat{W}_a - \hat{W}_c\right) - \tilde{W}_a^TW'(f + g\hat{u}).$$

Using Young's inequality, (12), (17), and (18) the Lyapunov derivative can be upper bounded as

$$\dot{V}_L \leq -\underline{q}\|\zeta\|^2 - \left(\underline{c} - \frac{k_a}{2}\right)\left\|\tilde{W}_c\right\|^2 - \frac{k_a}{2}\left\|\tilde{W}_a\right\|^2$$
$$+ \iota_c\left\|\tilde{W}_c\right\| + \iota_a\left\|\tilde{W}_a\right\| + \iota.$$

Completing the squares, the upper bound on the Lyapunov derivative may be written as

$$\dot{V}_L \leq -\frac{\underline{q}}{2}\|\zeta\|^2 - \left(\frac{\underline{c}}{2} - \frac{k_a}{4}\right)\left\|\tilde{W}_c\right\|^2 - \frac{k_a}{4}\left\|\tilde{W}_a\right\|^2$$
$$+ \frac{\iota_c^2}{2\underline{c} - k_a} + \frac{\iota_a^2}{2k_a} + \iota,$$

which can be further upper bounded as

$$\dot{V}_L \leq -\alpha\|Z\|, \, \forall\|Z\| \geq K > 0. \tag{24}$$

Using (19), (20), (21), and (24), Theorem 4.18 in [22] is invoked to conclude that $Z$ is ultimately bounded, in the sense that $\limsup_{t\to\infty}\|Z(t)\| \leq \underline{v_L}^{-1}(\overline{v_L}(K))$.

Based on the definition of $Z$ and the inequalities in (19) and (24), $\zeta, \tilde{W}_c, \tilde{W}_a \in \mathcal{L}_\infty$. Since $\zeta \in \mathcal{L}_\infty$ and $W$ is a continuous function of $\zeta$, $W(\zeta) \in \mathcal{L}_\infty$. Hence, $\hat{W}_c, \hat{W}_a \in \mathcal{L}_\infty$, which implies $\hat{u} \in \mathcal{L}_\infty$. From the definitions of $u_s$ and $s$, $\overline{u} \in \mathcal{L}_\infty$. ∎

## IV. SIMULATION RESULTS

Simulation results are provided to demonstrate the performance of the developed ADP-based path-planner. The simulation is performed for the control affine system given in (3), where $\overline{f}(\zeta) = \mathbf{0}$ and $\overline{g}(\zeta) = I_{2\times2}$.

For this particular example, the smooth scheduling function is defined as

$$s(\zeta) \triangleq \sum_{i=1}^{M}\begin{cases} 0.95, & \|\zeta - c_{obs_i}\| \leq r_{obs_i} \\ 0.95T(\zeta), & r_{obs_i} < \|\zeta - c_{obs_i}\| \leq r_{pen_i} \\ 0, & \text{otherwise} \end{cases} \tag{25}$$

where

$$T(\zeta) \triangleq \left(\frac{1}{2} + \frac{1}{2}\cos\left(\frac{\pi}{r_{pen_i} - r_{obs_i}}(\|\zeta - c_{obs_i}\| - r_{obs_i})\right)\right),$$

$M$ is the number of obstacles, $r_{obs_i}$ is a positive constant indicating the radius of the $i^{th}$ obstacle, $r_{pen_i}$ is a positive constant indicating the radius corresponding to the unsafe region surrounding the $i^{th}$ obstacle, and $c_{obs_i} \in \mathbb{R}^n$ denotes the center corresponding to the $i^{th}$ obstacle. With this formulation of the smooth scheduling function, it is assumed that the obstacles are selected such that unsafe regions do not overlap[2].

The continuous auxiliary controller $u_s$ is defined as

$$u_s(\zeta) \triangleq \frac{u_{sat}(\zeta - c_{obs_i})}{\|\zeta - c_{obs_i}\|}. \tag{26}$$

[2]If a group of obstacles are close enough for the unsafe regions to overlap, then the group may be considered as one large obstacle.
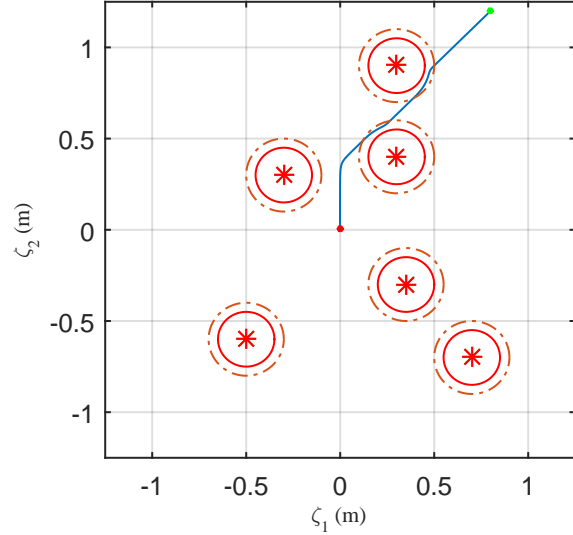


Figure 2. The path generated by the developed method is shown, where the dashed lines around the obstacles denote the boundary of the unsafe region and the solid lines denote the boundary of the obstacle.

In Appendix A, the auxiliary controller in (26) is shown to prevent the agent from entering the interior of an obstacle.

The unsafe region is where the penalty function $P$ begins to effect the agent's cost. The non-negative function $P$ is given as

$$P(\zeta) = \sum_{i=1}^{M}\begin{cases} 40 & \|\zeta - c_{obs_i}\| \leq r_{obs_i} \\ 40T(\zeta), & r_{obs_i} < \|\zeta - c_{obs_i}\| \leq r_{pen_i} \\ 0, & \text{otherwise} \end{cases}.$$

The initial state is selected as $\zeta_0 = \begin{bmatrix} 0.8 & 1.2 \end{bmatrix}^T$. The StaF basis for the value function approximation is selected as

$$\sigma = \begin{bmatrix} \zeta^T(\zeta + d_1) & \zeta^T(\zeta + d_2) & \zeta^T(\zeta + d_3) \end{bmatrix}^T,$$

where the centers of the kernels are selected as $d_1 = 0.005 \cdot \begin{bmatrix} 0 & 1 \end{bmatrix}^T$, $d_2 = 0.005 \cdot \begin{bmatrix} 0.8660 & -0.5 \end{bmatrix}^T$, $d_3 = 0.005 \cdot \begin{bmatrix} -0.8660 & -0.5 \end{bmatrix}^T$. The Bellman error is extrapolated to 25 sampled data points that are selected on a uniform $5 \times 5$ grid that spans a square of size 0.01, and is centered about the current state. The weighting matrices are selected as $\mathcal{Q} = I_{2\times2}$ and $\mathcal{R} = I_{2\times2}$. The learning gains are selected as $k_{c1} = 0.25$, $k_{c2} = 0.15$, $k_a = 0.5$, $\beta = 0.3$, and $k_\rho = 0.05$. The least squares update law's initial condition is selected as $\Gamma_0 = 300 \cdot I_{3\times3}$. The policy and value function weight estimates are arbitrarily initialized to

$$\hat{W}_c(0) = \hat{W}_a(0) = \begin{bmatrix} 0.4 & 0.4 & 0.4 \end{bmatrix}^T.$$

The generated path for both simulation trials are shown in Figure 2. Note that the state trajectories in Figure 2 briefly enter the unsafe region, where the auxiliary controller successfully escorts the agent away from the obstacle.

## V. Conclusion

An online approximation of a robust optimal path-planning strategy is developed. The solution to the HJB equation is approximated using adaptive dynamic programming. Since the unknown value function is approximated locally, the locations of the static obstacles do not need to be known until the obstacles are within an approximation window. The developed feedback policy guarantees ultimately bounded convergence of the approximated path to the optimal path without the requirement of persistent excitation, typically required for online adaptive dynamic programming. The results are validated with simulations. Future work will focus on comparisons with existing pregenerative literature to evaluate the accuracy of the developed method.

## Appendix A
## Auxiliary Controller Analysis

Consider a change of coordinates, where $\overline{\zeta} = \zeta - c_{obs_i}$. A positive definite function $V_{obs} : [0, \infty) \to \mathbb{R}^n$ is given as

$$V_{obs} = \overline{\zeta}^T \overline{\zeta}. \tag{27}$$

The time derivative of (27) is

$$\dot{V}_{obs} = 2\overline{\zeta}^T \dot{\zeta}.$$

Substituting the dynamics (1) with the definitions of $\overline{f}$ and $\overline{g}$ provided in Section IV, and the controller in (2) yields

$$\dot{V}_{obs} = 2\overline{\zeta}^T \left( s\left(\zeta\right) u_s + \left(1 - s\left(\zeta\right)\right) u\left(t\right) \right).$$

Substituting the auxiliary controller defined in (26) and the fact that the norm of the optimal controller $u$ is bounded by $\sqrt{2} u_{sat}$, the derivative is lower bounded by

$$\dot{V}_{obs} \geq 2 u_{sat} \|\zeta\| \left( s\left(\zeta\right) - \frac{\sqrt{2}}{1 + \sqrt{2}} \right). \tag{28}$$

Let $B_{obs_i}$ denote the local domain of the obstacle centered at $c_{obs_i}$ defined as $B_{obs_i} \triangleq \left\{ \zeta | \overline{\zeta} \leq r_{obs_i} \right\}$. By the definition of the scheduling function in (25),

$$\inf_{\zeta \in B_{obs_i}} s\left(\zeta\right) = 0.95. \tag{29}$$

Consider the inequality in (28) on the local domain $B_{obs_i}$, then (28) is further bounded by

$$\dot{V}_{obs} \geq 2 u_{sat} \|\zeta\| \left( \sqrt{2} - 1.05 \right).$$

Substituting the function $V_{obs}$, the derivative may be written as

$$\dot{V}_{obs} \geq 2 u_{sat} V_{obs}^{\frac{1}{2}} \left( \sqrt{2} - 1.05 \right).$$

Solving the differential equation using separation of variables, yields

$$V_{obs} \geq \left( \sqrt{2} - 1.05 \right)^2 u_{sat}^2 t^2$$

Hence, the obstacle center $c_{obs_i}$ is unstable in the local domain $B_{obs_i}$. Furthermore, a state trajectory starting outside the local domain $B_{obs_i}$ will not enter the interior of $B_{obs_i}$.

## References

[1] A. Alvarez, A. Caiti, and R. Onken, "Evolutionary path planning for autonomous underwater vehicles in a variable ocean," vol. 29, pp. 418–429, 2004.

[2] A. V. Rao, D. A. Benson, C. L. Darby, M. A. Patterson, C. Francolin, and G. T. Huntington, "Algorithm 902: GPOPS, A MATLAB software for solving multiple-phase optimal control problems using the Gauss pseudospectral method," *ACM Trans. Math. Softw.*, vol. 37, no. 2, pp. 1–39, 2010.

[3] K. Yang, S. K. Gan, and S. Sukkarieh, "An efficient path planning and control algorithm for RUAVs in unknown and cluttered environments," *J. Intell. Robot Syst.*, vol. 57, pp. 101–122, 2010.

[4] S. Karaman and E. Frazzoli, "Sampling-based algorithms for optimal motion planning," *Int. J. Rob. R*, vol. 30, pp. 846–894, 2011.

[5] S. M. LaValle and P. Konkimalla, "Algorithms for computing numerical optimal feedback motion strategies," *Int. J. Rob. Res.*, vol. 20, pp. 729–752, 2001.

[6] S. LaValle, *Planning Algorithms*. Cambridge University Press, 2006.

[7] A. Shum, K. Morris, and A. Khajepour, "Direction-dependent optimal path planning for autonomous vehicles," *Robot. and Auton. Syst.*, 2015.

[8] C. Petres, Y. Pailhas, P. Patron, Y. Petillot, J. Evans, and D. Lane, "Path planning for autonomous underwater vehicles," vol. 23, pp. 331–341, 2007.

[9] K. Vamvoudakis and F. Lewis, "Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem," *Automatica*, vol. 46, no. 5, pp. 878–888, 2010.

[10] D. Vrabie and F. Lewis, "Neural network approach to continuous-time direct adaptive optimal control for partially unknown nonlinear systems," *Neural Netw.*, vol. 22, no. 3, pp. 237 – 246, 2009.

[11] S. Bhasin, R. Kamalapurkar, M. Johnson, K. G. Vamvoudakis, F. L. Lewis, and W. E. Dixon, "A novel actor-critic-identifier architecture for approximate optimal control of uncertain nonlinear systems," *Automatica*, vol. 49, no. 1, pp. 89–92, 2013.

[12] H. Modares, F. Lewis, and M.-B. Naghibi-Sistani, "Adaptive optimal control of unknown constrained-input systems using policy iteration and neural networks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 24, no. 10, pp. 1513–1525, 2013.

[13] H. Modares, F. L. Lewis, and M.-B. Naghibi-Sistani, "Integral reinforcement learning and experience replay for adaptive optimal control of partially-unknown constrained-input continuous-time systems," *Automatica*, vol. 50, no. 1, pp. 193–202, 2014.

[14] R. Kamalapurkar, J. A. Rosenfeld, and W. E. Dixon, "State following (StaF) kernel functions for function approximation Part II: Adaptive dynamic programming," in *Proc. Am. Control Conf.*, 2015, pp. 521–526.

[15] M. Abu-Khalaf, F. L. Lewis, and J. Huang, "Hamilton-Jacobi-Isaacs formulation for constrained input nonlinear systems," in *Proc. IEEE Conf. Decis. Control*, vol. 5, 2004, pp. 5034–5040.

[16] S. Lyshevski, "Optimal control of nonlinear continuous-time systems: design of bounded controllers via generalized nonquadratic functionals," in *Proc. Am. Control Conf.*, 1998.

[17] H. Zhang, Y. Luo, and D. Liu, "Neural-network-based near-optimal control for a class of discrete-time affine nonlinear systems with control constraints," *IEEE Trans. Neural Netw.*, vol. 20, no. 9, pp. 1490–1503, 2009.

[18] R. Kamalapurkar, P. Walters, and W. E. Dixon, "Concurrent learning-based approximate optimal regulation," in *Proc. IEEE Conf. Decis. Control*, Florence, IT, Dec. 2013, pp. 6256–6261.

[19] P. Ioannou and J. Sun, *Robust Adaptive Control*. Prentice Hall, 1996.

[20] W. E. Dixon, A. Behal, D. M. Dawson, and S. Nagarkatti, *Nonlinear Control of Engineering Systems: A Lyapunov-Based Approach*. Birkhauser: Boston, 2003.

[21] M. Abu-Khalaf and F. Lewis, "Nearly optimal control laws for nonlinear systems with saturating actuators using a neural network HJB approach," *Automatica*, vol. 41, no. 5, pp. 779–791, 2005.

[22] H. K. Khalil, *Nonlinear Systems*, 3rd ed. Upper Saddle River, NJ, USA: Prentice Hall, 2002.