

# Online optimal control of nonlinear discrete-time systems using approximate dynamic programming

Travis DIERKS<sup>1</sup>, Sarangapani JAGANNATHAN<sup>2</sup>

1.DRS Sustainment Systems, Inc., 201 Evans Lane, St. Louis, MO 63121, U.S.A.;

2.Department of Electrical and Computer Engineering, Missouri University of Science and Technology, Rolla, MO 65409, U.S.A.

**Abstract:** In this paper, the optimal control of a class of general affine nonlinear discrete-time (DT) systems is undertaken by solving the Hamilton Jacobi-Bellman (HJB) equation online and forward in time. The proposed approach, referred normally as adaptive or approximate dynamic programming (ADP), uses online approximators (OLAs) to solve the infinite horizon optimal regulation and tracking control problems for affine nonlinear DT systems in the presence of unknown internal dynamics. Both the regulation and tracking controllers are designed using OLAs to obtain the optimal feedback control signal and its associated cost function. Additionally, the tracking controller design entails a feedforward portion that is derived and approximated using an additional OLA for steady state conditions. Novel update laws for tuning the unknown parameters of the OLAs online are derived. Lyapunov techniques are used to show that all signals are uniformly ultimately bounded and that the approximated control signals approach the optimal control inputs with small bounded error. In the absence of OLA reconstruction errors, an optimal control is demonstrated. Simulation results verify that all OLA parameter estimates remain bounded, and the proposed OLA-based optimal control scheme tunes itself to reduce the cost HJB equation.

**Keywords:** Online nonlinear optimal control; Hamilton Jacobi-Bellman; Online approximators; Discrete-time systems

## 1 Introduction

Online approximators (OLAs) have been widely used in the controller designs for uncertain discrete-time (DT) nonlinear systems [1]; however, system stability is typically the sole consideration although optimality is generally preferred. While linear systems accompanied by quadratic cost functions can achieve optimal control by solving the Riccati equation [2], the optimal control of nonlinear DT systems often requires solving the nonlinear Hamilton-Jacobi-Bellman (HJB) equation. To extend the results of linear optimal control theory to nonlinear systems, the state-dependent Riccati equation [3] is proposed for suboptimal control under certain tight assumptions including the need for full system dynamics.

To avoid solving the HJB equation, approximate solutions to the HJB equation have been proposed [4–6]. In [4], a Chebyshev series was proposed for approximating the system dynamics, boundary conditions, and cost function. In [5, 6], neural networks (NNs) are utilized to solve the DT nonlinear optimal regulation in an offline manner by ignoring NN reconstruction errors and assuming complete system dynamics. Recently, online methods to solve the continuous HJB equation were presented in [7] for linear systems using online policy iterations based on adaptive control.

In addition to the optimal regulation problem, the optimal tracking control problem has also recently been considered. In [8], the authors consider the  $H_\infty$  optimal tracking control

by linearizing the error equations about the origin yielding a locally optimal control law. The effort in [9] considers the receding horizon optimal tracking control by linearizing the nonlinear error dynamics about the origin. In [10], authors consider the HJB equation and employ similar techniques such as [6] to find an offline solution to the optimal tracking control problem.

In this work, a novel approach to the optimal regulation of nonlinear DT systems is adopted to solve the HJB equation online. Using an initial stabilizing control, an OLA is tuned online at each time step to learn the HJB equation in contrast to normal approaches to ADP [7], which use two indices: a time index for the system dynamics and a policy iteration index for the cost function approximator. Then, a second OLA is utilized that minimizes the HJB function based on the information provided by the first OLA. Knowledge of the internal system dynamics is not required while the control coefficient matrix alone is needed although it can be relaxed by introducing an additional OLA [11].

Subsequently, this approach is extended to include optimal tracking control even with unknown internal dynamics by using a third OLA to approximate the feedforward part of the control input that is normally required for tracking [1, 2]. Novel online parameter tuning laws for the OLAs are derived, and Lyapunov theory is utilized to demonstrate the stability of the system while explicitly considering the OLA approximation errors in contrast to the other works [5, 6, 10].

Received 29 July 2010; revised 11 April 2011.

This work was partly supported by the National Science Foundation (No.ECCS#0621924, ECCS-#0901562), and the Intelligent Systems Center.

© South China University of Technology and Academy of Mathematics and Systems Science, CAS and Springer-Verlag Berlin Heidelberg 2011

## 2 Background

Consider the affine nonlinear DT system

$$x_{k+1} = f(x_k) + g(x_k)u(x_k), \tag{1}$$

where  $x_k \in \mathbb{R}^n$ ,  $f(x_k) \in \mathbb{R}^n$ ,  $g(x_k) \in \mathbb{R}^{n \times m}$  satisfies  $\|g(x_k)\|_F \leq g_M$ , where  $\|\cdot\|_F$  denotes the Frobenius norm [1], and  $u(x_k) \in \mathbb{R}^m$  is the control input. Without loss of generality, assume that  $x = 0$  is a unique equilibrium point on a compact set  $\Omega$  while the states are considered measurable. In order to control (1) in an optimal manner, select the control sequence  $u(x_k)$  that minimizes the infinite horizon cost function as [6]

$$J(x_k) = Q(x_k) + u(x_k)^T R u(x_k) + J(f(x_k) + g(x_k)u(x_k)) \tag{2}$$

for all  $x_k$ , where with  $Q(x_k) > 0$  and  $R \in \mathbb{R}^{m \times m}$  is a symmetric positive definite matrix. Selecting  $Q(x_k) > 0$  ensures that variations in any direction of the state affect the cost, which can be linked to the observability condition [2]. In addition, it is required that the control policy  $u(x_k)$  be admissible [5] and  $J(x_k = 0) = 0$  so that  $J(x_k)$  serves as a Lyapunov function.

Using Bellman’s principle of optimality [2], the infinite horizon optimal cost function,  $J^*(x_k)$ , is time invariant and satisfies the DT HJB equation [2]

$$J^*(x_k) = \min_{u(x_k)} (r(x_k, u(x_k)) + J(f(x_k) + g(x_k)u(x_k))),$$

where  $r(x_k, u(x_k)) = Q(x_k) + u(x_k)^T R u(x_k)$ . The optimal control  $u^*(x_k)$  that minimizes  $J^*(x_k)$  is found by applying the stationary condition (see [2])

$$\frac{\partial J^*(x_k)}{\partial u(x_k)} = 2Ru(x_k) + g(x_k)^T \frac{\partial J^*(x_{k+1})}{\partial x_{k+1}} = 0,$$

and is shown to be (see [2])

$$u^*(x_k) = -\frac{1}{2}R^{-1}g(x_k)^T \frac{\partial J^*(x_{k+1})}{\partial x_{k+1}}. \tag{3}$$

The optimal control (3) is generally unavailable for nonlinear DT systems due to its dependence on the future state vector  $x_{k+1}$ . To circumvent this deficiency, a new approach to online optimal control is presented next.

## 3 Near optimal regulation of DT systems

The proposed approach entails two OLAs: one OLA to learn the HJB equation and a second OLA to learn the control signal that minimizes the estimated HJB equation. Using the function approximation property of OLAs [1], the cost function (2) and feedback control policy (3) have OLA representations on a compact set expressed as

$$J(x_k) = \Phi^T \sigma(x_k) + \varepsilon_{Jk}, \tag{4}$$

and

$$u(x_k) = \Theta^T \vartheta(x_k) + \varepsilon_{uk}, \tag{5}$$

respectively, where  $\Phi \in \mathbb{R}^{L_1}$  and  $\Theta \in \mathbb{R}^{L_2 \times m}$  are the constant target OLA parameters,  $\varepsilon_{Jk} \in \mathbb{R}$  and  $\varepsilon_{uk} \in \mathbb{R}^m$  are the bounded approximation errors, and  $\sigma(\cdot) \in \mathbb{R}^{L_1}$  and  $\vartheta(\cdot) \in \mathbb{R}^{L_2}$  are the activation functions for the cost and control signal OLA schemes, respectively. Here,  $L_1$  and  $L_2$  define the dimension of  $\Phi$  and  $\Theta$ , respectively. The upper bounds for the ideal OLA parameters are taken as  $\|\Phi\| \leq \Phi_M$  and  $\|\Theta\|_F \leq \Theta_M$ , where  $\Phi_M$  and  $\Theta_M$  are positive constants [1], respectively, and the approximation

errors are assumed to be bounded above as  $\|\varepsilon_{Jk}\| \leq \varepsilon_{JM}$  and  $\|\varepsilon_{uk}\| \leq \varepsilon_{uM}$ , where  $\varepsilon_{JM}$  and  $\varepsilon_{uM}$  are positive constants [1]. In addition, the gradient of the approximation error is assumed to be bounded above as  $\left\| \frac{\partial \varepsilon_{Jk}}{\partial x_{k+1}} \right\|_F \leq \varepsilon'_{JM}$ , where  $\varepsilon'_{JM}$  is also a positive constant [12].

To begin, the HJB function OLA is considered next.

### 3.1 Approximation of the optimal cost function

Now, the cost function (2) is approximated by an OLA and written as

$$\hat{J}(x_k) = \hat{\Phi}_k^T \sigma(x_k), \tag{6}$$

where  $\hat{J}(x_k)$  represents an approximated value of the cost function (2), and  $\hat{\Phi}_k$  is the estimate of the target OLA parameter vector  $\Phi$ . The basis function should satisfy  $\|\sigma(0)\| = 0$  for  $\|x\| = 0$  to ensure  $\hat{J}(0) = 0$  can be satisfied [2].

It is observed that (2) represents a nonlinear DT Lyapunov equation that can be rewritten as  $J(f(x_k) + g(x_k)u(x_k)) - J(x_k) + r(x_k, u(x_k)) = 0$ . However, this relationship is not guaranteed to hold when the estimated cost function in (6) is used. Therefore, the residual or cost-to-go (CTG) error associated with (6) can be written as

$$\hat{J}(x_{k+1}) - \hat{J}(x_k) + r(x_k, u(x_k)) = e_{Jk}$$

or

$$e_{Jk} = r(x_k, u(x_k)) + \hat{\Phi}_k^T \Delta \sigma(x_{k+1}), \tag{7}$$

where  $\Delta \sigma(x_{k+1}) = \sigma(x_{k+1}) - \sigma(x_k)$ . Next, we define an auxiliary CTG error vector as

$$E_{Jk} = Y_k + \hat{\Phi}_k^T X_k \in \mathbb{R}^{1 \times (1+j)}, \tag{8}$$

where

$$Y_k = [r(x_k, u(x_k)) \ r(x_{k-1}, u(x_{k-1})) \ \dots \ r(x_{k-j}, u(x_{k-j}))],$$

$$X_k = [\Delta \sigma(x_{k+1}) \ \Delta \sigma(x_k) \ \dots \ \Delta \sigma(x_{k+1-j})],$$

$0 < j < k - 1 \in \mathbb{N}$  and  $\mathbb{N}$  being the set of natural real numbers. It is clear that (8) represents a time history of the previous  $j + 1$  CTG errors (7) recalculated using the most recent  $\hat{\Phi}_k$ .

Moving on, define the OLA parameter update law as

$$\hat{\Phi}_{k+1} = \hat{\Phi}_k - \frac{\alpha_J X_k E_{Jk}^T}{\|X_k X_k^T + I\|_F}, \tag{9}$$

where  $0 < \alpha_J < 1$  is a small positive design parameter and  $I$  is an identity matrix of appropriate dimension. It is observed that the auxiliary CTG error (8) becomes zero when  $x_k = 0$  because the cost functions (4) and (6) are zero at  $x_k = 0$ . Thus, once the system states have converged to zero, the cost function OLA is no longer updated. This can be viewed as a persistency of excitation (PE) requirement for the inputs to the cost function OLA wherein the system, states must be persistently existing long enough for the optimal cost function to be learned.

To obtain the OLA parameter estimation error dynamics, rewrite (2) using the target OLA representation (6) as

$$r(x_k, u(x_k)) = -\Phi^T \Delta \sigma(x_{k+1}) - \Delta \varepsilon_{Jk}, \tag{10}$$

where  $\Delta \varepsilon_{Jk} = \varepsilon_{Jk+1} - \varepsilon_{Jk}$ . Substituting (10) into (7) results in

$$e_{Jk} = -\hat{\Phi}_k^T \Delta \sigma(x_{k+1}) - \Delta \varepsilon_{Jk}, \tag{11}$$

where  $\tilde{\Phi}_k = \Phi - \hat{\Phi}_k$  is the cost parameter estimation error. Similarly, (8) can be rewritten as

$$E_{Jk} = -\tilde{\Phi}_k^T X_k - \Psi_k, \tag{12}$$

where  $\Psi_k = [\Delta\varepsilon_{Jk} \ \Delta\varepsilon_{Jk-1} \ \cdots \ \Delta\varepsilon_{Jk-j}]$  and  $\|\Psi_k\|^2 \leq \Psi_M^2$ . Now, observing  $\tilde{\Phi}_{k+1} = \Phi - \hat{\Phi}_{k+1}$  and using (9) and (12) results in the OLA parameter estimation error dynamics to be expressed as

$$\tilde{\Phi}_{k+1} = \tilde{\Phi}_k + \frac{\alpha_J X_k E_{Jk}^T}{\|X_k X_k^T + I\|_F}. \tag{13}$$

Next, the action network that generates the optimal feedback signal (3) is considered.

### 3.2 Estimation of the optimal feedback control signal

The objective of this section is to find the control policy that minimizes the approximated cost function (6). To begin the development of the feedback control policy, we define the OLA approximation of (5) to be

$$\hat{u}(x_k) = \hat{\Theta}_k^T \vartheta(x_k), \tag{14}$$

where  $\hat{\Theta}_k$  is the estimated value of  $\Theta_k$ .

Next, the action error is defined as the difference between the feedback control applied to (1) and the control signal that minimizes the estimated cost function (6), which is denoted as

$$\tilde{u}(x_k) = \hat{\Theta}_k^T \vartheta(x_k) + R^{-1} g^T(x_k) \left( \frac{\partial \sigma(x_{k+1})}{\partial x_{k+1}} \right)^T \frac{\hat{\Phi}_k}{2}, \tag{15}$$

and the control OLA parameter update is defined to be

$$\hat{\Theta}_{k+1} = \hat{\Theta}_k - \frac{\alpha_u \vartheta(x_k) \tilde{u}_k^T}{(\vartheta^T(x_k) \vartheta(x_k) + 1)}, \tag{16}$$

where  $0 < \alpha_u < 1$  is a positive design parameter. Because the control policy  $u(x_k)$  in (5) minimizes the cost function (4), from (3) we can write

$$0 = \varepsilon_{uk} + \Theta^T \vartheta(x_k) + \frac{1}{2} R^{-1} g^T(x_k) \frac{\partial \varepsilon_{Jk+1}}{\partial x_{k+1}} + \frac{1}{2} R^{-1} g^T(x_k) \frac{\partial \sigma(x_{k+1})}{\partial x_{k+1}} \Phi. \tag{17}$$

Then, defining the action OLA parameter estimation error as  $\tilde{\Theta}_k = \Theta - \hat{\Theta}_k$  and subtracting (16) from (15) yields

$$\tilde{u}(x_k) = -\tilde{\Theta}_k^T \vartheta(x_k) - \frac{1}{2} R^{-1} g^T(x_k) \frac{\partial \sigma(x_{k+1})}{\partial x_{k+1}} \tilde{\Phi}_k - \tilde{\varepsilon}_{uk}, \tag{18}$$

where  $\tilde{\varepsilon}_{uk} = \varepsilon_{uk} + \frac{1}{2} R^{-1} g^T(x_k) \frac{\partial \varepsilon_{Jk+1}}{\partial x_{k+1}}$  and satisfies  $\|\tilde{\varepsilon}_{uk}\| \leq \tilde{\varepsilon}_{uM}$  for a positive constant  $\tilde{\varepsilon}_{uM}$ .

As a final step, we form the parameter estimation error dynamics as

$$\tilde{\Theta}_{k+1} = \tilde{\Theta}_k + \frac{\alpha_u \vartheta(x_k) \tilde{u}_k^T}{(\vartheta^T(x_k) \vartheta(x_k) + 1)}, \tag{19}$$

and the closed loop nonlinear system dynamics (1) can be rewritten in terms of  $u^*(x_k)$  and  $\hat{\Theta}_k$  as

$$\begin{aligned} x_{k+1} &= f(x_k) + g(x_k) \hat{u}(x_k) \\ &= f(x_k) + g(x_k) \hat{\Theta}_k^T \vartheta(x_k) \\ &= f(x_k) + g(x_k) u^*(x_k) - g(x_k) \tilde{\Theta}_k^T \vartheta(x_k) \\ &\quad - g(x_k) \varepsilon_{uk}. \end{aligned} \tag{20}$$

**Remark 1** To calculate the action error (14) and implement the OLA parameter update (15), knowledge of the

input transformation matrix  $g(x_k)$  is required. However, knowledge of the internal dynamics  $f(x_k)$  can be avoided by rewriting (2) in terms of  $x_{k+1}$  as noted in [6]. In contrast to [6], this work utilizes online learning and requires an initial stabilizing control to ensure the nonlinear system remains stable while the optimal control signal and cost function are computed because nonlinear systems are known to have finite escape times [13].

The following theorem will show that the cost function and action network OLA parameter estimation errors are uniformly ultimately bounded (UUB) [1]. Furthermore, the estimated control input (14) approaches the optimal control signal with small bounded error, which is a function of the OLA reconstruction errors  $\varepsilon_J$  and  $\varepsilon_u$ , and if the OLA approximation errors are considered to be negligible [5, 6], asymptotic convergence is observed.

The initial system states are considered to reside in a compact set  $\Omega \subset \mathbb{R}^n$  by using the initial stabilizing control input  $u_{0k}$ . Moreover, sufficient conditions for the OLA tuning gains,  $\alpha_J$  and  $\alpha_u$ , are derived to ensure that all future states never leave the compact set. As a result, in the compact  $\Omega$ , the cost function and its gradient as well as the basis function of the action network are bounded according to  $\|\sigma(x_k)\| \leq \sigma_M$ ,  $\left\| \frac{\partial \sigma(x_k)}{\partial x_k} \right\| \leq \sigma'_M$ , and  $\|\vartheta(x_k)\| \leq \vartheta_M$ , respectively. A similar approach has been used in several other well known efforts in control theory [14] to establish bounds on functions of the system states  $x_k$ .

**Theorem 1** (Convergence of the optimal control signal) Let  $u_0(x_k)$  be any initial admissible control policy for the class nonlinear controllable systems in (1). Let the OLA parameter tuning for the cost function estimator and the action network be provided by (9) and (15), respectively. Then, there exist positive constants  $\alpha_J$  and  $\alpha_u$  such that the system states and the cost and action network OLA parameter estimate errors  $\tilde{\Phi}_k$  and  $\tilde{\Theta}_k$ , respectively, are all UUB for all  $k \geq k_0 + T$ . Furthermore,  $\|\hat{u}(x_k) - u^*(x_k)\| \leq \varepsilon_r$  for a small positive constant  $\varepsilon_r$ .

Proof of Theorem 1 is shown in Appendix.

**Remark 2** The cost and control OLA parameter update laws are both written in terms of the time index  $k$  instead of a policy iteration index commonly used in standard ADP approaches [7]. Therefore, the conclusions of Theorem 1 illustrate the convergence of the ADP scheme and the boundedness of the system states simultaneously.

In the next section, the optimal online regulator is extended to consider the optimal tracking problem.

## 4 Near optimal tracking control

In this section, the optimal tracker is considered as an extension of regulation [2]. In addition to the assumptions on the system dynamics (1) presented in Section 2, in this section it is assumed that  $m = n$  so that  $g(x_k)$  is invertible. The objective for the optimal tracking problem is to find the optimal control sequence to make the nonlinear system in (1) track a desired trajectory  $x_{dk}$ . To begin the development, define the dynamics of the desired trajectory as (see [8])

$$x_{dk+1} = f(x_{dk}) + g(x_k) u_d(x_{dk}, x_k), \tag{21}$$

where  $f(x_{dk})$  is the internal dynamics of the nonlinear system (1) rewritten in terms of the desired state  $x_{dk}$ ,  $g(x_k)$  is the input transformation matrix presented in (1), and  $u_d(x_{dk})$  is the control input to the desired system. It is useful to note that when  $x_{dk}$ ,  $x_{dk+1}$ , and  $f(x_{dk})$  are known,  $u_d(x_{dk}, x_k)$  is obtained by rearranging (20) to give

$$u_d(x_{dk}, x_k) = g(x_k)^{-1}(x_{dk+1} - f_d(x_{dk})). \quad (22)$$

Next, define the tracking error as

$$e_k = x_k - x_{dk}. \quad (23)$$

Then, using (1) and (20), the tracking error dynamics are

$$\begin{aligned} e_{k+1} &= f(x_k) + g(x_k)u(x_k) - x_{dk+1} \\ &= f_e(e_k) + g(x_k)u_e(e_k), \end{aligned} \quad (24)$$

where  $f_e(e_k) = f(x_k) - f(x_{dk})$  and

$$u_e(e_k) = u(x_k) - u_d(x_{dk}, x_k). \quad (25)$$

Considering  $u_e(e_k)$  as the control input for (23),  $u_e(e_k)$  is an admissible control policy with  $e_k = 0$  being an equilibrium point of (23). To convert the nonlinear tracking into a regulation problem, the infinite horizon cost function (2) is rewritten in terms of  $e_k$  and  $u_e(e_k)$  as

$$J_e(e_k) = Q_e(e_k) + u_e(e_k)^T R_e u_e(e_k) + J_e(e_{k+1}), \quad (26)$$

with  $Q_e(e_k) > 0$  and  $R_e \in \mathbb{R}^{m \times m}$  is positive definite. Because  $u_e(e_k)$  is admissible, (26) is finite. The optimal control input that minimizes (26) is found by solving  $\frac{\partial J_e(e_k)}{\partial u_e(e_k)} = 0$  as

$$u_e^*(e_k) = -\frac{1}{2}R_e^{-1}g^T(x_k)\frac{\partial J_e^*(e_{k+1})}{\partial e_{k+1}}, \quad (27)$$

or

$$u^*(x_k) = u_d(x_{dk}) - \frac{1}{2}R_e^{-1}g^T(x_k)\frac{\partial J_e^*(e_{k+1})}{\partial e_{k+1}}. \quad (28)$$

It is observed that the optimal tracking control input (28) consists of a predetermined feedforward term,  $u_d(x_{dk}, x_k)$ , and an optimal feedback term that is a function of the gradient of the optimal cost function. Additionally, implementation of the feedforward term requires knowledge of the internal dynamics  $f(x_k)$  and control coefficient matrix  $g(x_k)$ .

In this effort,  $f(x_k)$ ,  $J_e(e_k)$  and  $\frac{\partial J_e(e_{k+1})}{\partial e_{k+1}}$  are all unknown. To mitigate these deficiencies, the universal approximation properties [1] of OLAs are utilized as described next.

The proposed solution for achieving nonlinear optimal tracking control entails three portions: an HJB function estimator that evaluates the performance of the error system by approximating (26), a feedback system that is designed to produce a nearly optimal portion of the control signal (27), and a feedforward design to produce the feedforward control input (28) by approximating  $f(x_{dk})$ . Using the approximation property of OLAs [1], the cost function (26), feedback control policy (27), and desired internal dynamics,  $f(x_{dk})$ , in (20) have OLA representations on a compact set expressed as

$$J_e(e_k) = \Phi_e^T \sigma_e(e_k) + \varepsilon_{J_e}, \quad (29)$$

$$u_e(e_k) = \Theta_e^T \vartheta_e(e_k) + \varepsilon_{u_e}, \quad (30)$$

and

$$f(x_{dk}) = \Omega_d^T \phi(x_{dk}) + \varepsilon_d(x_{dk}), \quad (31)$$

respectively, where  $\Phi_e$ ,  $\Theta_e$ , and  $\Omega_d$  are the constant target OLA parameters,  $\varepsilon_{J_e}$ ,  $\varepsilon_{u_e}$ , and  $\varepsilon_d$  are the bounded approximation errors, and  $\sigma_e(\cdot)$ ,  $\vartheta_e(\cdot)$  and  $\phi(\cdot)$  are the linearly independent vector activation functions for the cost, feedback, and feedforward control networks, respectively [1]. The upper bounds for the ideal OLA parameters are taken as  $\|\Phi_e\| \leq \Phi_{eM}$ ,  $\|\Theta_e\|_F \leq \Theta_{eM}$ , and  $\|\Omega_d\|_F \leq \Omega_{dM}$  where  $\Phi_{eM}$ ,  $\Theta_{eM}$ , and  $\Omega_{dM}$  are positive constants [1] while the approximation errors are assumed to be bounded above such that  $\|\varepsilon_{J_e}\| \leq \varepsilon_{J_eM}$ ,  $\|\varepsilon_{u_e}\| \leq \varepsilon_{u_eM}$ , and  $\|\varepsilon_d\| \leq \varepsilon_{dM}$  where  $\varepsilon_{J_eM}$ ,  $\varepsilon_{u_eM}$  and  $\varepsilon_{dM}$  are positive constants [1]. In addition, the gradient of the cost function approximation error is considered to be bounded according to  $\left\| \frac{\partial \varepsilon_{J_e}}{\partial e_{k+1}} \right\|_F \leq \varepsilon'_{J_eM}$

where  $\varepsilon'_{J_eM}$  is also a positive constant [12]. As in Section 3, the basis functions and the gradient of the cost function basis vector are considered to be bounded on a compact set [14] according to  $\|\sigma_e(\cdot)\| \leq \sigma_{eM}$ ,  $\|\vartheta_e(\cdot)\| \leq \vartheta_{eM}$ ,  $\|\phi(\cdot)\| \leq \phi_M$ , and  $\left\| \frac{\partial \sigma_e(\cdot)}{\partial(\cdot)} \right\| \leq \sigma'_{eM}$  for known constants  $\sigma_{eM}$ ,  $\vartheta_{eM}$ ,  $\phi_M$ , and  $\sigma'_{eM}$ , respectively [1]. To begin, the design of the HJB function and optimal feedback approximators will be considered first.

#### 4.1 Cost function and optimal feedback control

The objective of the optimal tracking control law is to stabilize system (23) while minimizing the cost function (26). To begin, (26) is approximated by an OLA as

$$\hat{J}_e(e_k) = \hat{\Phi}_{ek}^T \sigma_e(e_k), \quad (32)$$

where  $\hat{\Phi}_{ek}$  is the approximation for the ideal parameters  $\Phi_e$ . The basis vector  $\sigma_e(\cdot)$  is selected to satisfy  $\|\sigma_e(0)\| = 0$  to facilitate  $J_e(0) = 0$ .

Similar to Section 3, it is observed that (26) represents a nonlinear DT Lyapunov equation for tracking; however, when the estimated cost function in (32) is used, a residual or CTG error,  $e_{J_{ek}}$ , associated with (32) for tracking is required. The CTG error for tracking is defined identically to the CTG error for regulation, defined in (7) in Section 3, with  $x_k$ ,  $u(x_k)$  and  $\hat{\Phi}_k^T \Delta \sigma(x_{k+1})$  replaced by  $e_k$ ,  $u_e(e_k)$ , and  $\hat{\Phi}_{ek}^T \Delta \sigma_e(e_{k+1})$ , respectively. In addition, the auxiliary CTG error for tracking,  $E_{J_{ek}}$ , is defined similarly to (8) with the regulation variables replaced by those for tracking as described above.

The OLA parameter update law is now defined as

$$\hat{\Phi}_{ek+1} = \hat{\Phi}_{ek} - \frac{\alpha_{J_e} X_{ek} E_{J_{ek}}^T}{\|X_{ek} X_{ek}^T + I\|_F}, \quad (33)$$

where  $0 < \alpha_{J_e} < 1$  is a small positive design parameter,  $X_{ek}$  is defined similarly to  $X_k$  in (8), and  $I$  is an identity matrix of appropriate dimension. Following the steps described by (10)–(13) in Section 3, the OLA parameter estimation error dynamics for tracking  $\tilde{\Phi}_{ek+1}$  are written similarly to (13).

Next, the action network that generates the optimal feedback signal that minimizes the approximated cost function (32) is considered. To begin, define the OLA approximation of (30) as

$$\hat{u}_e(e_k) = \hat{\Theta}_{ek}^T \vartheta_e(e_k), \quad (34)$$

where  $\hat{\Theta}_{ek}$  is the OLA estimate of the target value  $\Theta_e$ . The

basis function,  $\vartheta_e(\cdot)$ , is selected to satisfy  $\|\vartheta_e(0)\| = 0$  so that  $\|\hat{u}_e(0)\| = 0$  can be satisfied as required for admissibility.

Next, the action error for tracking,  $\tilde{u}_e(e_k)$ , is defined to be the difference between the feedback control applied to (23) and the control signal that minimizes the estimated cost function (32), which is written similarly (14) in Section 3 with the regulation variables replaced by the appropriate variables defined for tracking. The proposed control OLA parameter update is defined to be

$$\hat{\Theta}_{ek+1} = \hat{\Theta}_{ek} - \frac{\alpha_{ue}\vartheta_e(e_k)\tilde{u}_{ek}^T}{(\vartheta_e^T(e_k)\vartheta(e_k) + 1)}, \quad (35)$$

where  $0 < \alpha_{ue} < 1$  is a positive design parameter. Defining the action OLA parameter estimation error for tracking as  $\tilde{\Theta}_{ek} = \Theta_e - \hat{\Theta}_{ek}$  and following similar steps as those presented in (16) and (17) in Section 3, the action error for tracking can be written identically to (17) for regulation, and the parameter estimation error dynamics  $\tilde{\Theta}_{ek+1}$ , for tracking can be written identically to (18).

Next, the overall nearly optimal control signal is derived and system stability is investigated.

#### 4.2 Nearly optimal control input

Recall that the optimal tracking control input (28) to system (1) is comprised of an optimal feedback term and a predetermined feedforward term, which is a function of the desired yet unknown internal dynamics,

$$f(x_{dk}) \equiv f(x_k)|_{x_k=x_{xd}}.$$

Therefore, the desired internal dynamics are approximated online by identifying the actual internal dynamics, and reevaluating the estimate using the desired state,  $x_{dk}$ .

To begin the online identifier development, the nonlinear system dynamics (1) are rewritten as

$$x_{k+1} = \Omega_d^T \phi(x_k) + \varepsilon_d(x_k) + g(x_k)u(x_k, x_{dk}, e_k), \quad (36)$$

where  $\Omega_d^T \phi(x_k) + \varepsilon_d(x_k)$  is the OLA estimate defined in (31) reevaluated at  $x_k$ . As in the previous cases, it is assumed that there exists a nonzero lower bound such that  $\phi_{dm} \leq \|\phi(\cdot)\|$  for a positive constant  $\phi_{dm}$ . The online identifier is then defined as

$$\hat{x}_{k+1} = \hat{\Omega}_{dk}^T \phi(x_k) + g(x_k)u(x_k, x_{dk}, e_k) - K_d \tilde{x}_k, \quad (37)$$

where  $\hat{x}_{k+1}$  is the estimate of  $x_{k+1}$ ,  $\hat{\Omega}_{dk}$  is the parameter estimate of  $\Omega_{dk}$ ,  $\tilde{x}_k = x_k - \hat{x}_k$ , and  $0 < K_d < 1$  is a constant. Next, subtracting (37) from (36) and defining  $\tilde{\Omega}_{dk} = \Omega_d - \hat{\Omega}_{dk}$ , the identifier error dynamics are

$$\tilde{x}_{k+1} = \tilde{\Omega}_{dk}^T \phi(x_k) + \varepsilon_d(x_k) + K_d \tilde{x}_k. \quad (38)$$

The parameter update law is now defined as

$$\hat{\Omega}_{dk+1} = \hat{\Omega}_{dk} + \alpha_d \phi(x_k)(\tilde{x}_{k+1} + K_d \tilde{x}_k)^T, \quad (39)$$

where  $0 < \alpha_d < 1$  is a constant, and the parameter estimation error dynamics,  $\tilde{\Omega}_{dk+1} = \Omega_d - \hat{\Omega}_{dk+1}$ , are found to be

$$\tilde{\Omega}_{dk+1} = \tilde{\Omega}_{dk} - \alpha_d \phi(x_k)(\tilde{\Omega}_{dk}^T \phi(x_k) + \varepsilon_d(x_k))^T. \quad (40)$$

Moving on, the feedforward control input can now be de-

finied using the OLA representation (31) as

$$u_d(x_{dk}, x_k) = g(x_k)^{-1}(x_{dk+1} - (\Omega_d^T \phi(x_{dk}) + \varepsilon_d(x_{dk}))). \quad (41)$$

Now, the estimated feedforward control input is defined as

$$\hat{u}_d(x_{dk}, x_k) = g(x_k)^{-1}(x_{dk+1} - \hat{\Omega}_{dk}^T \phi(x_{dk})), \quad (42)$$

where  $\hat{f}(x_{dk}) = \hat{\Omega}_{dk}^T \phi(x_{dk})$  was used while observing the definition of  $f(x_{dk})$  in (20). Now, using (34) and (42), the estimate of the control input (28) is written as

$$u(x_k, x_{dk}, e_k) = \hat{u}_d(x_{dk}, x_k) + \hat{u}_e(e_k), \quad (43)$$

and applying (43) to the nonlinear system (1) reveals

$$x_{k+1} = f(x_k) + g(x_k)(\hat{u}_d(x_{dk}, x_k) + \hat{u}_e(e_k))$$

or

$$e_{k+1} = f(x_k) - \hat{\Omega}_{dk}^T \phi(x_{dk}) + g(x_k)\hat{u}_e(e_k). \quad (44)$$

Then, adding and subtracting  $f(x_{dk})$  and  $g(x_k)u_e(e_k)$  to (44), recalling the OLA representations of  $f(x_{dk})$  and  $u_e(e_k)$  in (31) and (30), respectively, (44) is rewritten as

$$e_{k+1} = f_e(e_k) + g(x_k)u_e(e_k) + \tilde{\Omega}_{dk}^T \phi(x_{dk}) - g(x_k)\tilde{\Theta}_{ek}^T \vartheta(e_k) + \varepsilon_{ed}(x_{dk}), \quad (45)$$

where  $\varepsilon_{ed} = \varepsilon_d(x_{dk}) - g(x_k)\varepsilon_{ue}$ .

Next, the convergence of tracking error, online identifier, and the cost function, feedback control signal, and feedforward control signal OLA parameter estimation errors is demonstrated in the following theorem while explicitly considering the OLA reconstruction errors.

**Theorem 2** (System stability) Let  $u_{e0k}$  be any initial admissible control for the nonlinear system (23) such that (23) is initially asymptotically stable. Let the parameter tuning for the cost OLA and feedforward OLA be provided by (33) and (35), respectively, and let the tuning law for the feedforward estimator be given by (39). Then, there exist positive constants  $K_d$ ,  $\alpha_{Je}$ ,  $\alpha_{ue}$ , and  $\alpha_d$  such that the tracking error, identification error, and the OLA parameter estimation errors of the cost function, feedback and feedforward terms are all UUB for all  $k \geq k_0 + T$ , and the estimated control input approaches the optimal control input such that  $\|u(x_k, x_{dk}, e_k) - u^*(x_k, x_{dk}, e_k)\| \leq \varepsilon_u$  for a small positive constant  $\varepsilon_u$ .

Proof of Theorem 2 is shown in Appendix.

## 5 Simulation results

To demonstrate the effectiveness of the nearly optimal nonlinear controller, a tracking example is presented under two scenarios: when the internal dynamics are known and when they are not. When the internal dynamics are known, only the cost and feedback OLAs are required to regulate the tracking error to zero. Subsequently, the feedforward OLA is added when the internal dynamics are not known. For example, a differentially driven nonholonomic mobile robot is considered whose discretized nonlinear system is described by [15]

$$v_{k+1} = \begin{bmatrix} v_{Rk+1} \\ v_{Lk+1} \end{bmatrix} = \begin{bmatrix} v_{Rk} \\ v_{Lk} \end{bmatrix} + Tf(v_k) + TM^{-1}\tau, \quad (46)$$

where  $f(v_k) = -M^{-1}(V(v_k)v_k + F(v_k))$ ,  $\tau$  is the control torque,  $v_{Rk}$  and  $v_{Lk}$  are the velocities of the right and

left wheels of the robot, respectively, and  $T$  is the sampling time. Furthermore,  $M$  is the inertial matrix,  $V(v_k)$  is the nonlinear Coriolis forces matrix, and  $F(v_k)$  is the nonlinear friction vector. The robot parameters used in the simulation were chosen to be as those presented in [14] while the sampling time is taken as  $T = 0.01$  s. The objective of the mobile robot is to track a virtual reference cart, and the desired right and left wheel velocities,  $v_{kd} = [v_{dRk} \ v_{dLk}]^T$ , which are generated online [14]. For this test, the translational and angular velocities of the reference cart were taken as  $v_{rk} = 1$  m/s and  $\omega_{rk} = 0.5 \sin(0.2\pi kT)$  rad/s, respectively.

To implement the control scheme, two-layer NNs are considered consisting of one layer of randomly assigned constant weights,  $v_N$ , in the first layer and one layer of tunable weights,  $\Theta_N$ , in the second layer. It has been shown that by randomly selecting the input layer weights  $v_N$ , the activation function forms a stochastic basis, and thus the approximation property holds for all inputs in a compact set [1]. Additionally, 10 hidden layer neurons were selected for both the cost function and feedback control OLAs while 25 hidden layer neurons were selected to estimate the feedforward signal of the control input. The activation function of the cost function OLA was selected as hyperbolic tangent squared to obtain an even linearly independent basis function. Conversely, the basis function of the feedback action OLA was selected as hyperbolic tangent after taking the partial derivative. Finally, radial basis functions were selected as activation functions for the feedforward control estimator.

The initial stabilizing control law was selected as

$$u_{e0}(e_k) = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} v_{Rk} - v_{dRk} \\ v_{Lk} - v_{dLk} \end{bmatrix} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} e_{1k} \\ e_{2k} \end{bmatrix}. \tag{47}$$

The control gains for the OLA-based optimal control scheme were selected as  $K_d = 0.1$ ,  $\alpha_{Je} = 0.1$ ,  $\alpha_{ue} = 0.1$  and  $\alpha_d = 0.1$ , and all tunable NN weights were initialized to zero. The simulation was run for 10 s (1000 time steps), and for the first 6 s, probing noise with mean zero and variance 1.4 was added to the system to ensure the persistency of excitation condition holds (see Section 3).

The resulting robot trajectory when the dynamics are known is shown in Fig. 1, where it is observed that the robot converges to the path of the virtual cart and maintains the desired course for the remainder of the test. The time histories of both the cost function and the feedback control signal parameter estimates are shown in Fig. 2. Examining the figure, it is clear that the parameter estimates converge to constant values and remain bounded consistent with Theorem 1. The cost function and feedback action network errors are shown in Fig. 3 where it is clear that both errors initially incur large values but then converge to a small bounded value near the origin.

Next, knowledge of the internal dynamics is removed, and the feedforward estimator presented in Section 4 is added to the robot control law. The robot trajectories for

this case were observed to be similar to the trajectory in Fig. 1 while the final values of the cost and action OLA parameters as well as the values of the cost and action errors were observed to be similar to those presented in Figs. 2 and 3, respectively. In addition, the difference between the actual feedforward control term and the estimated feedforward term is shown in Fig. 4. Here, the estimation error is found to be small and bounded consistent with the theoretical results of Theorem 2.

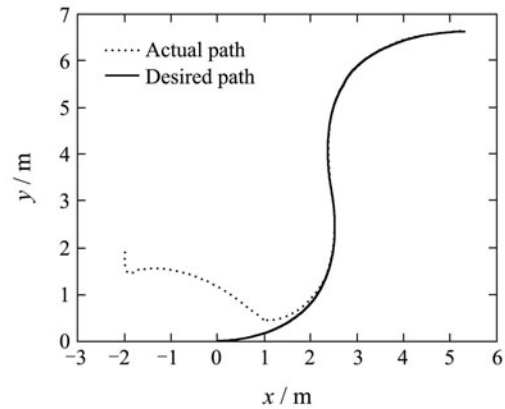


Fig. 1 Robot trajectory.

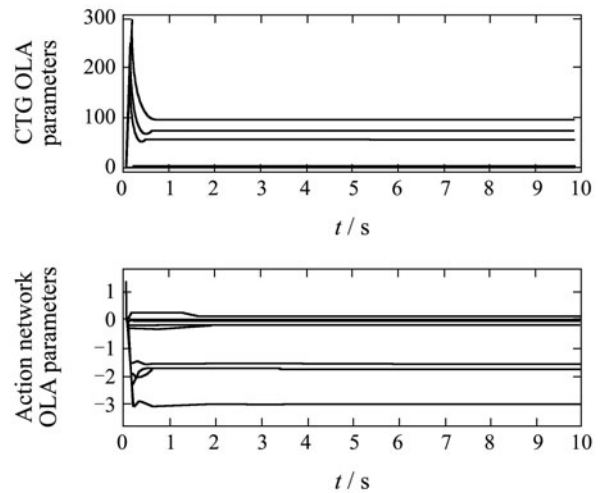


Fig. 2 OLA parameter estimates.

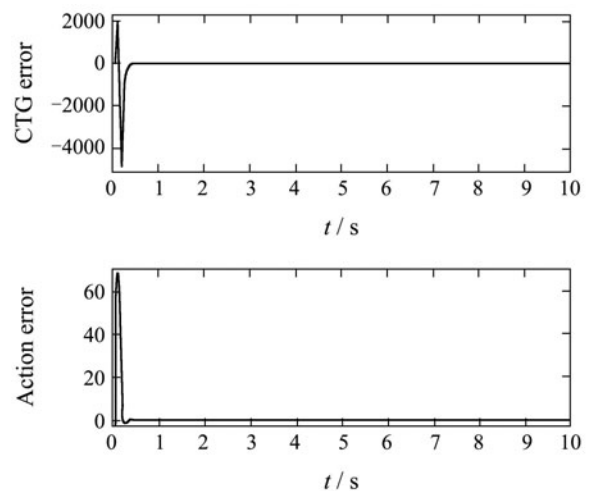


Fig. 3 OLA CTG and action errors.

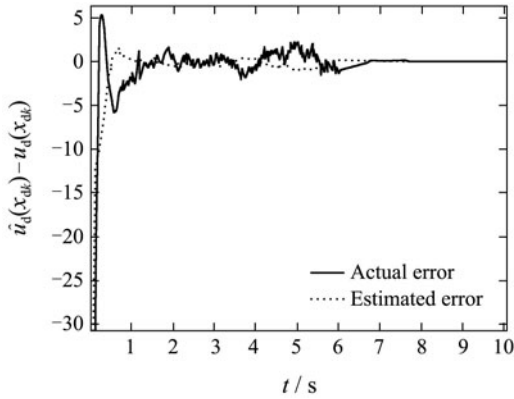


Fig. 4 OLA feedforward control term estimation error.

Closed-form solutions for the discrete-time HJB equation are difficult to obtain, and no benchmarking methods currently exist for evaluating discrete-time nonlinear optimal control laws. As a result, a cost comparison with the initial stabilizing control policy is considered where the costs were calculated by summing the cost function (26) from  $i = 0$  to  $i = 1000$ . Examining Table 1, the OLA-based optimal tracking control scheme is observed to incur less cumulative cost than the initial stabilizing control both when the internal dynamics were known and when they were unknown. Thus, the proposed tracker is indeed tuning to decrease the cost function (26) as predicted. Furthermore, the difference in the optimal costs when the dynamics are known and when they are not are observed to be small, which again reinforces the theoretical results of Theorem 2.

Table 1 Comparisons of total cost.

Control policy	Cost with $f_d$ known	Cost with $f_d$ unknown
$u_{e0}$	3.104e5	3.320e5
$\hat{u}_e^*$	2.562e5	2.565e5

## 6 Conclusions

In this work, the HJB equation was solved online for the nearly optimal control of general affine nonlinear DT systems using OLAs to address the regulation and tracking control problems. Knowledge of the system's internal dynamics was not needed while the OLAs generate a novel nearly optimal control law, and an initial admissible control policy guarantees that the system is stable while the OLAs learn. For the tracking problem, a desired feedforward portion of the control input generated by the third OLA rendered an overall stable system with ultimately bounded signals. All OLA parameters were tuned online using novel update laws, and system stability is guaranteed using Lyapunov theory.

## References

- [1] S. Jagannathan. *Neural Network Control of Nonlinear Discrete-time Systems*. Boca Raton: CRC Press, 2006.
- [2] F. L. Lewis, V. L. Syrmos. *Optimal Control*. 2nd ed. New York: John Wiley & Sons, 1995.
- [3] J. Shamma, J. Cloutier. Existence of SDRE stabilizing feedback. *IEEE Transactions on Automatic Control*, 2003, 48(3): 513 – 517.
- [4] J. Vlassenbroeck, R. Van Dooren. A Chebyshev technique for solving nonlinear optimal control problems. *IEEE Transactions on Automatic Control*, 1988, 33(4): 333 – 340.
- [5] Z. Chen, S. Jagannathan. Generalized Hamilton-Jacobi-Bellman

formulation based neural network control of affine nonlinear discrete-time systems. *IEEE Transactions on Neural Networks*, 2008, 19(1): 90 – 106.

- [6] A. Al-Tamimi, F. L. Lewis, M. Abu-Khalaf. Discrete-time nonlinear HJB solution using approximate dynamic programming: convergence proof. *IEEE Transactions on Systems, Man, and Cybernetics – Part B*, 2008, 38(4): 943 – 949.
- [7] D. Vrabie, O. Pastravanu, M. Abu-Khalaf, et al. Adaptive optimal control for continuous-time linear systems based on policy iteration. *Automatica*, 2009, 45(2): 477 – 484.
- [8] G. Toussaint, T. Basar, F. Bullo. H-infinity-optimal tracking control techniques for nonlinear underactuated systems. *Proceedings of IEEE Conference on Decision and Control*, New York: IEEE, 2000: 2078 – 2083.
- [9] D. Gu, H. Hu. Receding horizon tracking control of wheeled mobile robots. *IEEE Transactions on Control Systems Technology*, 2006, 14(4): 743 – 49.
- [10] H. Zhang, Q. Wei, Y. Luo. A novel infinite-time optimal tracking control scheme for a class of discrete-time nonlinear systems via the greedy HDP iteration algorithm. *IEEE Transactions on Systems, Man, and Cybernetics – Part B*, 2008, 38(4): 937 – 942.
- [11] T. Dierks, B. T. Thumati, S. Jagannathan. Optimal control of unknown affine nonlinear discrete-time systems using offline-trained neural networks with proof of convergence. *Neural Networks*, 2009, 22(5/6): 851 – 860.
- [12] K. G. Vamvoudakis, F. L. Lewis. Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem. *Automatica*, 2010, 46(5): 878 – 888.
- [13] H. K. Khalil. *Nonlinear Systems*, 3rd ed. New York: Prentice Hall, 2002.
- [14] D. Wang, J. Huang. Neural network-based adaptive dynamic surface control for a class of uncertain nonlinear systems in strict-feedback form. *IEEE Transactions on Neural Networks*, 2005, 16(1): 195 – 202.
- [15] M. K. Bugeja, S. G. Fabri, L. Camilleri. Dual adaptive dynamic control of mobile robots using neural networks. *IEEE Transactions on Systems, Man, and Cybernetics – Part B*, 2009, 39(1): 129 – 141.

## Appendix

**Proof of Theorem 1** Consider the positive definite Lyapunov candidate

$$V = \frac{\alpha_u^2 \alpha_J \bar{X}_m}{2g_M^2 \Omega_A (\vartheta_M^2 + 1)} V_D(x_k) + \frac{\alpha_J \bar{X}_m}{\Pi_A} V_u(\tilde{\Theta}_k) + V_J(\tilde{\Phi}_k), \quad (a1)$$

where

$$V_D(x_k) = x_k^T x_k, \quad V_u(\tilde{\Theta}_k) = \text{tr}\{\tilde{\Theta}_k^T \tilde{\Theta}_k\}, \quad V_J(\tilde{\Phi}_k) = \tilde{\Phi}_k^T \tilde{\Phi}_k, \\ \Pi_A = (\alpha_u + 1) (\sigma_M' g_M \lambda_{\max}(R^{-1}))^2 \\ + \alpha_u \frac{((\sigma_M' g_M \lambda_{\max}(R^{-1}))^2 + 2)}{4},$$

and  $\lambda_{\max}(R^{-1})$  is the maximum singular value of  $R$ . The first difference of (a1) is given by

$$\Delta V = \frac{\alpha_u^2 \alpha_J \bar{X}_m \Delta V_D(x_k)}{2g_M^2 \Pi_A (\vartheta_M^2 + 1)} + \frac{\alpha_J \bar{X}_m \Delta V_u(\tilde{\Theta}_k)}{\Pi_A} + \Delta V_J(\tilde{\Phi}_k). \quad (a2)$$

Considering first  $\Delta V_D(x_k) = x_{k+1}^T x_{k+1} - x_k^T x_k$ , substituting the nonlinear dynamics (19), and applying the Cauchy-Schwartz inequality  $((a_1 + a_2 + \dots + a_n)^2 \leq n(a_1^2 + a_2^2 + \dots + a_n^2))$  twice, reveals

$$\Delta V_D(x_k) \leq 2\|f(x_k) + g(x_k)u^*(x_k)\|^2 - x_k^T x_k \\ + 4\|g(x_k)\tilde{\Theta}^T \vartheta(x_k)\|^2 + 4\|g(x_k)\varepsilon_{uk}\|^2. \quad (a3)$$

It is observed that the optimal closed loop system,  $f(x_k) + g(x_k)u^*(x_k)$ , is asymptotically stable on a compact set. Furthermore, it can be shown that on a compact set, the optimal closed

loop system is upper bounded according to

$$\|f(x_k) + g(x_k)u^*(x_k)\|^2 \leq k^* \|x_k\|^2,$$

where  $k^*$  is a constant. Using this observation while defining  $\Xi_{Ak} = \tilde{\Theta}_k^T \vartheta(x_k)$ ,  $\Delta V_D(x_k)$  is upper bounded according to

$$\Delta V_D(x_k) \leq -(1 - 2k^*) \|x_k\|^2 + 4g_M^2 \|\Xi_{Ak}\|^2 + 2g_M^2 \varepsilon_{uM}^2. \tag{a4}$$

Moving on,  $\Delta V_u(\tilde{\Theta}_k) = \text{tr}\{\tilde{\Theta}_{k+1}^T \tilde{\Theta}_{k+1}\} - \text{tr}\{\tilde{\Theta}_k^T \tilde{\Theta}_k\}$ , and substituting (18) and then (17),  $\Delta V_u(\tilde{\Theta}_k)$  is upper bounded in terms of  $\Xi_{Ak}$  as

$$\begin{aligned} \Delta V_u(\tilde{\Theta}_k) &\leq -\frac{\alpha_u(2 - \alpha_u)}{(\vartheta^T(x_k)\vartheta(x_k) + 1)} \|\Xi_{Ak}\|^2 \\ &\quad + \frac{\alpha_u(\alpha_u + 1) \left( \|\tilde{\Phi}_k\| \sigma'_M g_M \lambda_{\max}(R^{-1}) + 2\varepsilon_{uM} \right)}{(\vartheta^T(x_k)\vartheta(x_k) + 1)} \|\Xi_{Ak}\| \\ &\quad + \frac{\alpha_u^2}{(\vartheta^T(x_k)\vartheta(x_k) + 1)} \cdot \left( \frac{(\sigma'_M g_M \lambda_{\max}(R^{-1}))^2 \|\tilde{\Phi}_k\|^2}{4} \right. \\ &\quad \left. + \sigma'_M g_M \lambda_{\max}(R^{-1}) \|\tilde{\Phi}_k\| \varepsilon_{uM} \right) + \alpha_u^2 \varepsilon_{uM}^2. \end{aligned}$$

Now, observing

$$-\alpha_u(2 - \alpha_u) \|\Xi_{Ak}\|^2 = -\frac{\alpha_u(3 - 3\alpha_u)}{2} \|\Xi_{Ak}\|^2 - \frac{\alpha_u(\alpha_u + 1)}{2} \|\Xi_{Ak}\|^2$$

and completing the square with respect to

$$\|\Xi_{Ak}\| \left( \|\tilde{\Phi}_k\| \sigma'_M g_M \lambda_{\max}(R^{-1}) + 2\varepsilon_{uM} \right),$$

and then  $\|\tilde{\Phi}_k\| \sigma'_M g_M \lambda_{\max}(R^{-1}) \varepsilon_{uM}$  reveals the upper bound of  $\Delta V_u(\tilde{\Theta}_k)$  to be written as

$$\begin{aligned} \Delta V_u(\tilde{\Theta}_k) &\leq -\frac{\alpha_u}{2(\vartheta_M^2 + 1)} (3 - 3\alpha_u) \|\Xi_{Ak}\|^2 + \alpha_u \Pi_A \|\tilde{\Phi}_k\|^2 \\ &\quad + \left( \alpha_u(5\alpha_u + 4) + \frac{\alpha_u^2 (\sigma'_M g_M \lambda_{\max}(R^{-1}))^2}{2} \right) \varepsilon_{uM}^2. \end{aligned} \tag{a5}$$

Next, considering

$$\Delta V_J(\tilde{\Phi}_k) = \tilde{\Phi}_{k+1}^T \tilde{\Phi}_{k+1} - \tilde{\Phi}_k^T \tilde{\Phi}_k,$$

and substituting (12) and (13) reveals

$$\begin{aligned} \Delta V_J(\tilde{\Phi}_k) &= -2\alpha_J \frac{\tilde{\Phi}_k^T X_k X_k^T \tilde{\Phi}_k}{\|X_k X_k^T + I\|_F} - \alpha_J \frac{\tilde{\Phi}_k^T X_k \Psi_k^T}{\|X_k X_k^T + I\|_F} \\ &\quad - \alpha_J \frac{\Psi_k X_k^T \tilde{\Phi}_k}{\|X_k X_k^T + I\|_F} + \alpha_J^2 \frac{\tilde{\Phi}_k^T X_k X_k^T X_k X_k^T \tilde{\Phi}_k}{\|X_k X_k^T + I\|_F^2} \\ &\quad + \alpha_J^2 \frac{\tilde{\Phi}_k^T X_k X_k^T X_k \Psi_k^T}{\|X_k X_k^T + I\|_F^2} + \alpha_J^2 \frac{\Psi_k X_k^T X_k X_k^T \tilde{\Phi}_k}{\|X_k X_k^T + I\|_F^2} \\ &\quad + \alpha_J^2 \frac{\Psi_k X_k^T X_k \Psi_k^T}{\|X_k X_k^T + I\|_F^2}. \end{aligned}$$

Observing

$$0 < \bar{X}_m \leq \frac{\|X_k^T X_k\|}{\|X_k X_k^T + I\|_F} < 1,$$

where  $\bar{X}_m$  is a constant ensured to exist by the PE condition described in Section 3,  $\Delta V_J(\tilde{\Phi}_k)$  can be upper bounded as

$$\Delta V_J(\tilde{\Phi}_k) \leq -\alpha_J \bar{X}_m \|\tilde{\Phi}_k\|^2 + \alpha_J^2 \Psi_M^2. \tag{a6}$$

Then, substituting (a4)–(a6) into (a2) yields the upper bound for  $\Delta V$  to be

$$\begin{aligned} \Delta V &\leq \frac{-(1 - 2k^*)\alpha_u^2 \alpha_J \bar{X}_m}{2g_M^2 \Pi_A (\vartheta_M^2 + 1)} \|x_k\|^2 - \alpha_J \bar{X}_m (1 - \alpha_u) \|\tilde{\Phi}_k\|^2 \\ &\quad - \frac{\alpha_u \alpha_J \bar{X}_m (3 - 7\alpha_u)}{2(\vartheta_M^2 + 1) \Pi_A} \|\Xi_{Ak}\|^2 + \varepsilon_{SM}, \end{aligned} \tag{a7}$$

where

$$\begin{aligned} \varepsilon_{SM} &= \alpha_J^2 \Psi_M^2 + \frac{\varepsilon_{uM}^2 \alpha_u^2 \alpha_J \bar{X}_m}{(\Pi_A (\vartheta_M^2 + 1))} \\ &\quad + \alpha_J \bar{X}_m \left( \alpha_u(5\alpha_u + 4) + \frac{\alpha_u^2 (\sigma'_M g_M \lambda_{\max}(R^{-1}))^2}{2} \right) \frac{\varepsilon_{uM}^2}{\Pi_A}. \end{aligned}$$

Therefore,  $\Delta V$  is less than zero provided the following inequalities hold

$$\begin{cases} \|x_k\| > \sqrt{\frac{2g_M^2 \Pi_A (\vartheta_M^2 + 1) \varepsilon_{SM}}{(1 - 2k^*) \alpha_u^2 \alpha_J \bar{X}_m}} \text{ or} \\ \|\Xi_{Ak}\| > \sqrt{\frac{\varepsilon_{SM} 2(\vartheta_M^2 + 1) \Pi_A}{\alpha_u \alpha_J \bar{X}_m (3 - 7\alpha_u)}} \equiv b_{\Xi A} \text{ or} \\ \|\tilde{\Phi}_k\| > \sqrt{\frac{\varepsilon_{SM}}{\alpha_J \bar{X}_m (1 - \alpha_u)}}, \end{cases} \tag{a8}$$

and the tuning gains are selected as  $\alpha_u < \frac{3}{7}$  and  $0 < \alpha_J < 1$  for the class of nonlinear systems that satisfy the optimal closed loop bounds described above with  $0 < k^* < \frac{1}{2}$ . Thus, using standard Lyapunov extensions [1], the system states and the cost and control NN weight estimation errors are UUB, and the system states are guaranteed to never leave their initial compact set.

To show  $\|\hat{u}(x_k) - u^*(x_k)\| \leq \varepsilon_r$ , use (5) and (14) to observe  $\hat{u}(x_k) - u^*(x_k) = -\tilde{\Theta}_k^T \vartheta(x_k) - \varepsilon_u$ . Then, taking the limit as  $k \rightarrow \infty$  and taking the upper bound of  $\hat{u}(x_k) - u^*(x_k)$  shows  $\|\hat{u}(x_k) - u^*(x_k)\| \leq \|\Xi_{Ak}\| + \varepsilon_{uM} \leq b_{\Xi A} + \varepsilon_{uM} \equiv \varepsilon_r$ , where  $b_{\Xi A}$  is defined in (a8).

**Proof of Theorem 2** Consider the positive definite Lyapunov candidate

$$\begin{aligned} V_T &= \frac{\alpha_{ue} \alpha_{Je} \alpha_d \bar{X}_{em}}{12g_M^2 \Pi_{Ae} (\vartheta_{eM}^2 + 1)} V_e(e_k) + \alpha_d V_{Je}(\tilde{\Phi}_{ek}) \\ &\quad + \frac{\alpha_{Je} \alpha_d \bar{X}_{em}}{\Pi_{Ae}} V_{eu}(\tilde{\Theta}_{ek}) \\ &\quad + \frac{\alpha_{Je} \alpha_{ue} \bar{X}_{em}}{2g_M^2 \Pi_{Ae} (\vartheta_{eM}^2 + 1)} V_d(\tilde{\Omega}_{dk}, \tilde{x}_k), \end{aligned} \tag{a9}$$

where

$$V_e(e_k) = e_k^T e_k, \quad V_{Je}(\tilde{\Phi}_{ek}) = \tilde{\Phi}_{ek}^T \tilde{\Phi}_{ek},$$

$$V_{eu}(\tilde{\Theta}_{ek}) = \text{tr}\{\tilde{\Theta}_{ek}^T \tilde{\Theta}_{ek}\},$$

$$V_d(\tilde{\Omega}_{dk}, \tilde{x}_k) = \text{tr}\{\tilde{\Omega}_{dk}^T \tilde{\Omega}_{dk}\} + \frac{\alpha_d^2 \tilde{x}_k^T \tilde{x}_k}{6},$$

and

$$\begin{aligned} \Pi_{Ae} &= (\alpha_{ue} + 1) (\sigma'_{eM} g_M \lambda_{\max}(R_e^{-1}))^2 \\ &\quad + \frac{\alpha_{ue} ((\sigma'_{eM} g_M \lambda_{\max}(R_e^{-1}))^2 + 2)}{4}, \end{aligned}$$

where  $\lambda_{\max}(R_e^{-1})$  is the maximum eigenvalue of  $R_e^{-1}$ . The first difference of  $V_T$  is given by

$$\begin{aligned} \Delta V_T &= \frac{\alpha_{ue} \alpha_{Je} \alpha_d \bar{X}_{em} \Delta V_e(e_k)}{12g_M^2 \Pi_{Ae} (\vartheta_{eM}^2 + 1)} + \alpha_d \Delta V_{Je}(\tilde{\Phi}_{ek}) \\ &\quad + \frac{\alpha_{Je} \alpha_d \bar{X}_{em}}{\Pi_{Ae}} \Delta V_{eu}(\tilde{\Theta}_{ek}) \\ &\quad + \frac{\alpha_{Je} \alpha_{ue} \bar{X}_{em} \Delta V_d(\tilde{\Omega}_{dk}, \tilde{x}_k)}{2g_M^2 \Pi_{Ae} (\vartheta_{eM}^2 + 1)}. \end{aligned} \tag{a10}$$

First, considering  $\Delta V_e(e_k)$ , substituting the error dynamics (45), and applying the Cauchy-Schwartz inequality twice yields

$$\begin{aligned} \Delta V_e(e_k) &\leq 2\|f_e(e_k) + g(x_k)u_e(e_k)\|^2 + 6\|\tilde{\Omega}_{dk}^T \phi(x_d)\|^2 \\ &\quad + 6g_M^2 \|\tilde{\Theta}_{ek}^T \vartheta(e_k)\|^2 + 6\varepsilon_{dM}^2 - e_k^T e_k. \end{aligned} \tag{a11}$$

Observing the optimal error dynamics are asymptotically stable, it can be shown that  $\|f_e(e_k) + g(x_k)u_e(e_k)\|^2 \leq k_e^* \|e_k\|^2$  where  $k_e^*$  is a constant. Applying this relation and taking the upper



bound of (a11) is written as

$$\Delta V_e(e_k) \leq -(1 - 2k_e^*) \|e_k\|^2 + 6 \|\tilde{\Omega}_{dk}^T \phi(x_d)\|^2 + 6g_M^2 \|\tilde{\Theta}_{ek}^T \vartheta(e_k)\|^2 + 6\varepsilon_{dM}^2. \quad (a12)$$

Next, examining the similarities of  $V_{Je}(\tilde{\Phi}_{ek})$  and  $V_{eu}(\tilde{\Theta}_{ek})$  to  $V_J(\tilde{\Phi}_k)$  and  $V_u(\tilde{\Theta}_k)$ , respectively, defined in Proof of Theorem 1, the first differences of  $\Delta V_{Je}(\tilde{\Phi}_{ek})$  and  $\Delta V_{eu}(\tilde{\Theta}_{ek})$  are calculated identically to (a5) and (a6), respectively, after substituting the appropriate variables defined for tracking in Section 4 in place of the regulation variables defined in Section 3.

Now, considering  $\Delta V_d(\tilde{\Omega}_{dk}, \tilde{x}_k)$ , using (38) and (40) and taking the upper bound gives

$$\begin{aligned} \Delta V_d(\tilde{\Omega}_{dk}, \tilde{x}_k) &\leq -\alpha_d(2 - \alpha_d) \|\tilde{\Omega}_{dk}^T \phi(x_k)\|^2 \\ &\quad + \frac{\alpha_d^2 \|\tilde{\Omega}_{dk}^T \phi(x_k)\|^2}{2} + 2\alpha_d(1 + \alpha_d) \|\tilde{\Omega}_{dk}^T \phi(x_k)\| \varepsilon_{dM} \\ &\quad + \varepsilon_{dM}^2 + \frac{\alpha_d^2 \varepsilon_{dM}^2}{2} + \frac{\alpha_d^2 K_d^2 \|\tilde{x}_k\|^2}{2} - \frac{\alpha_d^2 \|\tilde{x}_k\|^2}{6}. \end{aligned}$$

Then, completing with respect to  $\|\tilde{\Omega}_{dk}^T \phi(x_k)\| \varepsilon_{dM}$  yields

$$\begin{aligned} \Delta V_d(\tilde{\Omega}_{dk}, \tilde{x}_k) &\leq -\alpha_d \left(\frac{3}{2} - 2\alpha_d\right) \|\tilde{\Omega}_{dk}^T \phi(x_k)\|^2 \\ &\quad - \left(\frac{\alpha_d^2}{2}\right) \left(\frac{1}{3} - K_d^2\right) \|\tilde{x}_k\|^2 + (1 + 2\alpha_d(\alpha_d + 1) + \frac{\alpha_d^2}{2}) \varepsilon_{dM}^2. \end{aligned} \quad (a13)$$

Next, using (a12), (a13), and  $V_{Je}(\tilde{\Phi}_{ek})$  and  $V_{eu}(\tilde{\Theta}_{ek})$  written in the form of (a5) and (a6), respectively, to form (a10) gives

$$\begin{aligned} \Delta V_T &= -\frac{\alpha_{ue}\alpha_{Je}\alpha_d\bar{X}_{em}(1 - 2k_e^*)}{12g_M^2\Pi_{Ae}(\vartheta_{eM}^2 + 1)} \|e_k\|^2 \\ &\quad - \alpha_d\alpha_{Je}\bar{X}_{em}(1 - \alpha_{ue}) \|\tilde{\Phi}_{ek}\|^2 \\ &\quad - \frac{\alpha_{Je}\alpha_{ue}\alpha_d^2\bar{X}_{em}(\frac{1}{3} - K_d^2)}{4g_M^2\Pi_{Ae}(\vartheta_{eM}^2 + 1)} \|\tilde{x}_k\|^2 \\ &\quad - \frac{\alpha_{Je}\alpha_{ue}\alpha_d\bar{X}_{em}}{2g_M^2\Pi_{Ae}(\vartheta_{eM}^2 + 1)} \left(\frac{1}{2} - 2\alpha_d\right) \|\tilde{\Omega}_{dk}^T \phi(x_k)\|^2 \\ &\quad - \frac{\alpha_{ue}\alpha_{Je}\alpha_d\bar{X}_{em}}{2(\vartheta_{eM}^2 + 1)\Pi_{Ae}} (2 - 3\alpha_{eu}) \|\tilde{\Theta}_{ek}^T \vartheta(e_k)\|^2 + \varepsilon_{Se}, \end{aligned}$$

where

$$\begin{aligned} \varepsilon_{Se} &\equiv \frac{\varepsilon_{dM}^2\alpha_{Je}\alpha_{ue}\bar{X}_{em}}{2g_M^2\Pi_{Ae}(\vartheta_{eM}^2 + 1)} \left(1 + 2\alpha_d\left(\frac{5}{4}\alpha_d + \frac{3}{2}\right)\right) \\ &\quad + \alpha_d\alpha_{Je}^2\Psi_{eM}^2 + \frac{\alpha_{Je}\alpha_d\bar{X}_{em}}{\Pi_{Ae}} (\alpha_{ue}(5\alpha_{ue} + 4)) \\ &\quad + \frac{\alpha_{ue}^2(\sigma'_{eM}g_M\lambda_{\max}(R_e^{-1}))^2}{2} \varepsilon_{ueM}^2, \end{aligned}$$

and  $\Delta V_T < 0$  for the class of nonlinear systems that satisfy the optimal closed loop bounds described above with  $k_e^* < \frac{1}{2}$  provided the control gains are selected according to

$$K_d^2 < \frac{1}{3}, \quad \alpha_{Je} < 1, \quad \alpha_{eu} < \frac{2}{3}, \quad \alpha_d < \frac{1}{4},$$

and the following inequalities hold:

$$\begin{aligned} \|e_k\| &> \sqrt{\frac{12g_M^2\Pi_{Ae}(\vartheta_{eM}^2 + 1)\varepsilon_{Se}}{\alpha_{ue}\alpha_{Je}\alpha_d\bar{X}_{em}(1 - 2k_e^*)}} \equiv b_e \quad \text{or} \\ \|\tilde{x}_k\| &> \sqrt{\frac{4g_M^2\Pi_{Ae}(\vartheta_{eM}^2 + 1)\varepsilon_{Se}}{\alpha_{Je}\alpha_{ue}\alpha_d^2\bar{X}_{em}(\frac{1}{3} - K_d^2)}} \equiv b_{\tilde{x}} \quad \text{or} \\ \|\tilde{\Theta}_{ek}^T \vartheta(e_k)\| &> \sqrt{\frac{2(\vartheta_{eM}^2 + 1)\Pi_{Ae}\varepsilon_{Se}}{\alpha_{Je}\alpha_d\alpha_{ue}\bar{X}_{em}(2 - 3\alpha_{eu})}} \equiv b_{\Xi_e} \quad \text{or} \\ \|\tilde{\Phi}_{ek}\| &> \sqrt{\frac{\varepsilon_{Se}}{\alpha_d\alpha_{Je}\bar{X}_{em}(1 - \alpha_{ue})}} \equiv b_{\Phi_e} \quad \text{or} \\ \|\tilde{\Omega}_{dk}^T \phi(x_k)\| &> \sqrt{\frac{2g_M^2\Pi_{Ae}(\vartheta_{eM}^2 + 1)\varepsilon_{Se}}{\alpha_d\alpha_{Je}\alpha_{ue}\bar{X}_{em}(1/2 - 2\alpha_d)}} \equiv b_{\Omega_d}. \end{aligned}$$

Therefore, using standard Lyapunov extension [1], it can be concluded that  $\Delta V_T$  is less than zero outside of a compact set so that the tracking error (22) and the OLA parameter estimation errors of the cost function, feedback control signal, and feedforward control inputs are all UUB. To show  $\|u(x_k, x_{dk}, e_k) - u^*(x_k, x_{dk}, e_k)\| \leq \varepsilon_u$ , we use (28), (30), (34), and (41)–(43) to observe

$$\begin{aligned} &\|u(x_k, x_{dk}, e_k) - u^*(x_k, x_{dk}, e_k)\| \\ &= g_M^I \|\tilde{\Omega}_{kd}^T \phi(x_{dk})\|_F + g_M^I \varepsilon_{dM} + \|\tilde{\Theta}_{ek}^T \vartheta(e_k)\| + \varepsilon_{ueM} \\ &\leq g_M^I b_{\Omega_d} + g_M^I \varepsilon_{dM} + b_{\Xi_e} + \varepsilon_{ueM} \\ &\equiv \varepsilon_u. \end{aligned}$$



**Travis DIERKS** received his B.S. and M.S. degrees in Electrical Engineering and Ph.D. from the Missouri University of Science and Technology (formerly the University of Missouri-Rolla), Rolla, in 2005, 2007, and 2009, respectively. While completing his doctoral degree, he was a GAANN fellow and a Chancellor's fellow. His research interests include nonlinear control, optimal control, neural network control, and the control and coordination of autonomous ground and aerial vehicles. He is currently with DRS Sustainment Systems, Inc., St. Louis, MO. E-mail: tdierks@drs-ssi.com.



**Sarangapani JAGANNATHAN** received his Ph.D. degree in Electrical Engineering from the University of Texas, Arlington in 1994. He worked at Systems and Controls Research Division, in Caterpillar Inc., Peoria as a consultant during 1994 to 1998, and during 1998 to 2001 he was at the University of Texas at San Antonio as an assistant professor, and since September 2001, he is at the Missouri University of Science and Technology (former University of Missouri-Rolla) where he is currently a Rutledge-Emerson Distinguished Professor and Site Director for the NSF Industry/University Cooperative Research Center on Intelligent Maintenance Systems. He has coauthored over 87 peer reviewed journal articles, 175 refereed IEEE conference articles, several book chapters and three books entitled 'Neural Network Control of Robot Manipulators and Nonlinear Systems', published by Taylor & Francis, London in 1999, 'Discrete-time Neural Network Control of Nonlinear Discrete-time Systems', CRC Press, April 2006, and 'Wireless Ad Hoc and Sensor Networks: Performance, Protocols and Control', CRC Press, April 2007, and holds 18 patents. His research interests include adaptive and neural network control, computer/communication/sensor networks, prognostics, and autonomous systems/robotics.