



# Dimension reduction based adaptive dynamic programming for optimal control of discrete-time nonlinear control-affine systems

Qiang Li and Yunjun Xu

Department of Mechanical and Aerospace Engineering, University of Central Florida, Orlando, FL, USA

## ABSTRACT

Dynamic programming based methods have been widely used in solving discrete-time nonlinear constrained optimal control problems. However, applying these methods in real-time is challenging because a large amount of memory is needed and the associated computational cost is high. Here, a search space dimension reduction strategy is proposed for a class of nonlinear discrete-time systems that are control-affine and invertible. Specifically, a bio-inspired motion rule is combined with inverse dynamics to reduce the value iteration search space to one dimension. The corresponding suboptimal control algorithm is developed and its optimality is analysed. An adaptation rule is developed to estimate uncertainties and improve the base policy. The closed-loop system is proven to be asymptotically stable. The advantages of the algorithm including much smaller computational cost and significantly reduced memory usage are demonstrated with two simulation examples.

## ARTICLE HISTORY

Received 13 September 2021  
Accepted 2 August 2022

## KEYWORDS

Nonlinear constrained optimal control; adaptive dynamic programming; real-time control; control-affine system; uncertainty

## 1. Introduction

Nonlinear constrained optimal control aims to drive a nonlinear system towards a desired value or along a desired trajectory subject to equality and/or inequality constraints, while optimising a user-defined performance index (Betts, 2010; Pontryagin, 1986). This type of problem is ubiquitous and has been seen in many applications such as spacecraft maneuvering (Xin & Pan, 2009), wind turbine optimal control (Gros & Schild, 2017), and autonomous scouting in agricultural fields (N. Li et al., 2015).

Dynamic programming (DP), as one subset of nonlinear constrained optimal control methods, has been widely studied for many decades, especially for those with discrete-time dynamics (Chisci et al., 1998; Howard, 1960; Werbos, 1989a). Based on Bellman's Principle of Optimality, as shown in the seminal work by Bellman (1954), the basic procedure of DP involves two steps of operations over a discretised state-action search space. A value iteration (VI) algorithm is an explicit implementation of such a procedure (Bellman, 1961). In the first step, the Bellman Equation (Howard, 1960) is followed to compute the value function backward in time over the entire state-action space until the optimal value function is reached. In the second step, the optimal value function generates the greedy control actions along the state trajectory (Bellman, 1961). As another implementation of the procedure, in a policy iteration (PI) algorithm (Bellman, 1961), instead of searching all possible control actions, a pre-selected control policy is used in the first step to calculate the value function, which is then used to improve the control policy until an optimal control policy is reached. In the second step, the optimal control policy is followed along the state trajectory (Bellman, 1961). Memory usage

and computational cost of DP, especially in Step 1 of both VI and PI, grow exponentially with respect to the dimension of search space and the number of discretised nodes in each dimension – referred to as the curses of dimensionality (Bellman, 1961). This has prevented DP from being utilised in many real-time optimal control problems. For stochastic systems, this problem gets even worse.

For the purpose of mitigating the computational burden of DP, parameterised function approximators (Busoniu et al., 2010; Lendaris et al., 2002; Powell, 2007), primarily neural networks (Bertsekas & Tsitsiklis, 1996; Werbos, 1989b), have been applied to predict the values associated with unvisited states based on those of visited ones. However, the function approximation approaches have their downsides. To start with, it has been pointed out that replacing the lookup table in DP with function approximators sometimes is not robust even for very simple problems (Boyan & Moore, 1995). Besides, there is no guideline about the choice of function approximators, for example the configurations in neural network (Tesauro, 1992). A successful design of a function approximator highly relies on experiences and tuning (Powell, 2007). Additionally, a large amount of data and offline training are typically required to produce a close approximation of the value function (Polak, 1973). Furthermore, to obtain a good approximation, the number of parameters in a function approximator has to be sufficiently large and a good portion of the state-action space needs to be visited (Borghese & Arbib, 1995).

Another set of popular methods to mitigate the computational burden of DP are rollout algorithms (RA) (Bertsekas et al., 1997) and receding horizon control (RHC) (Mayne & Michalsha, 1988) or model predictive control (Clarke

et al., 1987) algorithms. Although RHC is not a DP exclusive method, it and RA utilise similar multi-step lookahead strategies and approximate the original optimal control problem with one of a shorter horizon at each time instant (Bertsekas, 2020). They essentially solve the Bellman equation in a neighbourhood of the current state. From this perspective, these methods can be regarded as subspace based methods that trade optimality for computational cost reductions. The difference between RA and RHC lies in the fact that the former approximates the value function outside the subspace following a base policy (Bertsekas, 2013), whereas the latter truncates the value function beyond the horizon (Michalsha & Mayne, 1993). Depending on the size of the subspace, which is determined by the number of lookahead steps, computational cost can still be a concern for real-time implementation of these methods (Bertsekas, 2013). On the other hand, techniques for reducing computational cost of standard DP can be implemented within the subspace of RA and RHC to further reduce the computational cost.

With the introduction of approximate/adaptive dynamic programming (ADP) (Barto et al., 1983; Prokhorov & Wunsch, 1997; Sutton, 1988; Watkins, 1989; Werbos, 1982), the offline computation of DP is transferred to online and the feedback information is then used to adapt to system uncertainties. ADP is a class of methods (Barto et al., 1983; Prokhorov & Wunsch, 1997; Sutton, 1988; Watkins, 1989; Werbos, 1982) that solve the Bellman equation forward in time through successive steps of policy evaluation and policy improvement. Among ADP methods, heuristic dynamic programming (HDP) (Werbos, 1990) sees the most applications (Bhatnagar, 2010; Borkar, 2005; Liu et al., 2013). It uses an actor network as the control policy for the online forward propagation and a critic network approximating the value function based on system feedback information (Werbos, 1990). The critic network output is then used to evaluate the actor network and improve it accordingly (Werbos, 1990). This procedure is repeated until an optimal control policy is reached (Werbos, 1990). As improvement of the control policy usually requires the calculations of value function derivatives, dual heuristic programming (DHP) is developed to approximate the derivative directly (Werbos, 1992). Furthermore, with globalised DHP (GDHP), not only the value function but also its derivative are approximated (Miller et al., 1995). There are also action-dependent (AD) variants of the aforementioned methods (i.e. ADHDP Watkins & Dayan, 1992, ADDHP Prokhorov & Wunsch, 1997, and ADGDHP Prokhorov & Wunsch, 1997) that assume the value function to be directly related to both state and action. However, because most policy-based ADP methods (Pi et al., 2020; Silver et al., 2014) assume the value function to be differentiable, they are not suitable for constrained optimal control problems where smoothness of the value function sometimes cannot be guaranteed (Yaghmaie & Braun, 2019). Additionally, although the computation of the value function is transferred from offline to online, a sufficient collection of online data is still necessary for the training of the critic network (Luo et al., 2018), resulting in high CPU occupation and memory usage.

In brief, it is still a challenging task to apply DP and its variations to many discrete-time, nonlinear constrained optimal control problems in real-time. Two observations are: (1) the inherent dynamics, many of which are known or partially

known, haven't been fully utilised; and (2) knowledge of the environment is not fully exploited, leading to inefficient search over noncontributory regions of the state-action space. A subspace method has been discussed to reduce the computational cost by finding local optimal or suboptimal solutions in Wei et al. (2017).

A two-step search space dimension reduction strategy is presented here to tackle the aforementioned issue for a class of nonlinear discrete-time dynamic systems that are control-affine and invertible. Instead of searching the action space exhaustively, the first dimension reduction step uses an inverse dynamics (ID) policy (Forrest-Barlach & Babcock, 1987; Kumar & Seywald, 1996) as the base policy to calculate the control action based on the current state and the desired state-to-reach in an Euler scheme, resulting in an ID-based VI (ID-VI) method. Furthermore, the state search space is confined using the virtual motion camouflage (VMC) rule, in which the state variation is constrained by the 1-dimensional path control parameter (PCP) (Xu & Basset, 2012). The motion camouflage phenomenon (Srinivasan & Davey, 1995) is observed in mating hoverflies when a male hoverfly is approaching a female one while its motion is shown static on the retina of the female counterpart. Xu and Basset (2012) have utilised this strategy in developing direct collocation based suboptimal controllers for continuous-time optimal trajectory planning problems, which are dramatically different from the problems and approaches in this paper. Here, the VMC rule is used in the second dimension reduction step of DP to project the state space into a 1-dimensional PCP parameter space and the ID-VI method is then developed into the VMC-based ID-VI (ID-VMC-VI) method. The main structure of VI method remains intact. Similar to most VI methods, once the value function has converged, the actual control action will be generated by the greedy policy (Powell, 2007).

The advantages of the resulting ID-VMC-VI algorithm are summarised as follows. (1) Since the search space is 1-dimensional, the optimisation related computation time is low. (2) Since only values corresponding to the 1-dimensional PCP variables are stored, the memory usage is significantly reduced. This can benefit many real-world applications with limited random access memory. (3) It is shown that the closed-loop system is asymptotically stable and the solution optimality is lower bounded by a first-order tracking controller or a second order tracking controller under two different scenarios. (4) When combined with a gradient based ID policy improvement rule, the ID-VMC-VI algorithm can solve constrained optimal control problems with uncertainties. A stability analysis is provided for the overall algorithm.

This study is a significant extension from our previous conference publication (Li & Xu, 2020). (1) Here we focus on the development of ID-VMC-VI method as a general method for discrete-time invertible nonlinear control-affine systems. In comparison, Li and Xu (2020) only applied the preliminary version of ID-VI, ID-VMC-VI, and another algorithm in an unmanned aerial vehicle (UAV) angular velocity control problem. (2) The theoretical aspect of ID-VMC-VI has been revised and extended with respect to the a general class of discrete-time optimal control problems, considering state and control inequality constraints in a systematic manner. (3) Additional sections on dimensionality, stability and optimality analyses

have been developed. (4) A different ID policy improvement method is developed to fit a wider range of problems with different performance indices. Additional convergence and stability analyses are provided. (5) Guidelines on how to apply the algorithm are offered. (6) In this paper, two simulation examples are used for illustration. While the first example is adopted and refined from Li and Xu (2020), the improved algorithm and new ID policy improvement method are applied. Furthermore, the simulation parameters and settings are tuned for a better result demonstration. The second example is a new and challenging application involving obstacle avoidance.

Section 2 formulates the problem and recalls necessary background knowledge on DP. The two search dimension reduction steps are presented in Section 3, along with stability, dimensionality and optimality analyses. Section 4 presents the ID policy improvement rule and its stability and convergence analyses. The overall algorithm is summarised in Section 5 and two simulation examples are shown in Section 6. The study is concluded at the end.

## 2. Problem definition and background information

A discrete-time, invertible, and control-affine system can be formulated as

$$\mathbf{x}_{t_{i+1}} = F(\mathbf{x}_{t_i}, \boldsymbol{\theta}) + G(\mathbf{x}_{t_i}, \boldsymbol{\theta})\mathbf{u}_{t_i}, \quad \mathbf{x}_{t_i} \in \mathcal{S}^n, \mathbf{u}_{t_i} \in \mathcal{A}^m, \boldsymbol{\theta} \in \mathcal{P}^r \quad (1)$$

where  $\mathbf{x}_{t_i}$  is the state vector at time instant  $t_i$ ,  $i = 0, 1, 2, \dots$ .  $\mathbf{u}_{t_i}$  is the control vector.  $\boldsymbol{\theta}$  consists of uncertain parameters.  $\mathcal{S}^n \subset \mathbb{R}^n$  is the constrained state vector space.  $\mathcal{A}^m \subset \mathbb{R}^m$  is the constrained control vector space.  $\mathcal{P}^r \subset \mathbb{R}^r$  is the uncertain parameter vector space. Assuming that  $G(\mathbf{x}_{t_i}, \boldsymbol{\theta})$  has a rank of  $m$ , its Moore-Penrose inverse  $G^{-1}(\mathbf{x}_{t_i}, \boldsymbol{\theta})$  exists (Penrose, 1954).

Because of the uncertainties in  $\boldsymbol{\theta}$ , given state  $\mathbf{x}_{t_i}$  and control  $\mathbf{u}_{t_i}$ , the predicted system state  $\hat{\mathbf{x}}_{t_{i+1}}$  based on  $\hat{\boldsymbol{\theta}}_{t_i}$ , which is the estimate of  $\boldsymbol{\theta}$  at  $t_i$ , will differ from the actual state feedback  $\mathbf{x}_{t_{i+1}}$ . Define  $\tilde{\mathbf{x}}_{t_{i+1}} = \mathbf{x}_{t_{i+1}} - \hat{\mathbf{x}}_{t_{i+1}}$  and let  $e_{t_{i+1}} = \tilde{\mathbf{x}}_{t_{i+1}}^T \tilde{\mathbf{x}}_{t_{i+1}} / 2$ . The estimation problem (E1) is defined as: Solve for  $\hat{\boldsymbol{\theta}}_{t_{i+1}}$  to minimise  $e_{t_{i+1}}$ .

The control problem (C1) is defined as: Solve for the optimal control sequence  $\{\mathbf{u}^*\}_{k=i}^{N-1}$  so that the following performance index, or value function, is minimised (Howard, 1960):

$$V(\mathbf{x}_{t_i}) = \sum_{k=i}^{N-1} L(\mathbf{x}_{t_k}, \mathbf{u}_{t_k}) \quad (2)$$

subject to the equality constraint (1) and the boundary conditions:

$$\mathbf{x}_{t_i} = \mathbf{x}_{t_i} \quad \text{and} \quad \mathbf{x}_{t_N} = \mathbf{x}_{d,t_i} \quad (3)$$

where  $L(\mathbf{x}_{t_k}, \mathbf{u}_{t_k}) \geq 0$  is the step cost of taking  $\mathbf{u}_{t_k}$  in  $\mathbf{x}_{t_k}$ .  $\mathbf{x}_{d,t_i} \in \mathcal{S}^n$  is the desired state-to-reach from  $\mathbf{x}_{t_i}$ .  $L(\mathbf{x}_{t_k}, \mathbf{u}_{t_k}) = 0$  only when  $\mathbf{x}_{t_k} = \mathbf{x}_{d,t_i}$  and  $F(\mathbf{x}_{d,t_i}, \hat{\boldsymbol{\theta}}_{t_i}) + G(\mathbf{x}_{d,t_i}, \hat{\boldsymbol{\theta}}_{t_i})\mathbf{u}_{t_k} = \mathbf{x}_{d,t_i}$ .

The basic VI scheme solves for the optimal value function  $V^*(\mathbf{x}_{t_i})$  in C1 over the  $(m+n)$ -dimensional state-action space

using the Bellman equation (Powell, 2007)

$$V_{j+1}(\mathbf{x}_{t_k}) = \min_{\mathbf{u}_{t_k} \in \mathcal{A}^m} \{L(\mathbf{x}_{t_k}, \mathbf{u}_{t_k}) + V_j(\mathbf{x}_{t_{k+1}})\}, \quad \forall \mathbf{x}_{t_k} \in \mathcal{S}^n \quad (4)$$

where  $k \geq i$ . With (4),  $V_j(\mathbf{x}_{d,t_i}) = 0$ ,  $j = 0, 1, 2, \dots$ , is guaranteed. State constraints are introduced in the form of penalty in the value function, i.e.  $V_j(\mathbf{x}_{t_k}) = \text{Penalty}$ ,  $\forall \mathbf{x}_{t_k} \notin \mathcal{S}^n$ .

In this paper, we consider the tabular approximation of the value function. We define  $\mathcal{S}^n$  to always refer to the discrete state vector space, and  $\mathcal{A}^m$  to refer to the discrete action vector space.

For an actual system, the discrete vector spaces  $\mathcal{S}^n$  and  $\mathcal{A}^m$  are finite. Therefore, within finite iterations, the value iteration stopping condition,  $V_{j+1}(\mathbf{x}_{t_k}) = V_j(\mathbf{x}_{t_k})$ ,  $\forall \mathbf{x}_{t_k} \in \mathcal{S}$ , will be met and the optimal value function  $V^*(\mathbf{x}_{t_k})$ ,  $\forall \mathbf{x}_{t_k} \in \mathcal{S}$ , will be reached. The optimal value function at  $t_i$  satisfies

$$V^*(\mathbf{x}_{t_k}) = \min_{\mathbf{u}_{t_k} \in \mathcal{A}^m} \{L(\mathbf{x}_{t_k}, \mathbf{u}_{t_k}) + V^*(\mathbf{x}_{t_{k+1}})\}, \quad \forall \mathbf{x}_{t_k} \in \mathcal{S}^n, \quad (5)$$

and the optimal control to take in state  $\mathbf{x}_{t_i}$  with respect to (2) is given as Powell (2007)

$$\mathbf{u}_{t_i}^* = \arg \min_{\mathbf{u}_{t_i} \in \mathcal{A}^m} \{L(\mathbf{x}_{t_i}, \mathbf{u}_{t_i}) + V^*(\mathbf{x}_{t_{i+1}})\}. \quad (6)$$

Discussions on the existence and uniqueness of  $\mathbf{u}_{t_i}^*$  can be found in Powell (2007). In general, it is assumed that 'arg min' includes all necessary operations to find the best control vector.

Because  $\mathbf{x}_{d,t_i}$  and  $\hat{\boldsymbol{\theta}}_{t_i}$  are time-varying, C1 needs to be solved at each time instant  $t_i$ . However, due to the curse of dimensionality, solving C1 via basic VI scheme can be impracticable when the system dimension is high, even if  $\mathcal{S}^n$  and  $\mathcal{A}^m$  are discrete and finite. We present a two-step search space dimension reduction strategy in the following.

## 3. Search space dimension reduction

### 3.1 Action space elimination

In the first dimension reduction step, the control action is calculated by an ID policy (Forrest-Barlach & Babcock, 1987) based on the current state and a state-to-reach. Following this procedure, instead of searching the  $(m+n)$ -dimensional state-action space to find the optimal control action, we only need to search the  $m$ -dimensional state space for the optimal state-to-reach. The optimal control action can then be calculated by the ID policy.

The ID policy based control to take in  $\mathbf{x}_{t_k}$  to reach  $\mathbf{x}_{d,t_k}$  in one time step is calculated based on the estimate  $\hat{\boldsymbol{\theta}}_{t_i}$  as

$$\begin{aligned} \mathbf{u}_{\pi,t_k} &= G^{-1}(\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_i})[\mathbf{x}_{d,t_k} - F(\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_i})] \\ &\triangleq \pi(\mathbf{x}_{t_k}, \mathbf{x}_{d,t_k}, \hat{\boldsymbol{\theta}}_{t_i}) \end{aligned} \quad (7)$$

where  $\mathbf{x}_{d,t_k}$  is assumed reachable within one time step without considering the constraints.  $\mathbf{u}_{\pi,t_k}$  is then saturated to satisfy the control constraint  $\mathbf{u}_{\pi,t_k} \in \mathcal{A}^m$  before used for the state propagation in the value iteration.

Choosing  $\mathbf{x}_{d,t_k} = \mathbf{x}_{d,t_i}$  and substituting (7) into (4), the Bellman equation for ID-VI can be written as

$$V_{\pi,j}(\mathbf{x}_{t_k}) = L(\mathbf{x}_{t_k}, \pi(\mathbf{x}_{t_k}, \mathbf{x}_{d,t_i}, \hat{\theta}_{t_i})) + V_{\pi,j-1}(\mathbf{x}_{t_{k+1}}), \quad \forall \mathbf{x}_{t_k} \in \mathcal{S}^n \quad (8)$$

where we follow the notation from Sutton and Barto (2018) and use  $V_{\pi}$  to denote the value function of following policy  $\mathbf{u}_{t_k} = \pi(\mathbf{x}_{t_k}, \mathbf{x}_{d,t_i}, \hat{\theta}_{t_i})$  during value iteration.

Compared with (4), which needs to be solved in the  $(m+n)$ -dimensional state-action space, (8) is solved in the  $n$ -dimensional state space. Thus, the action space is eliminated from the search space following ID-VI and the search space dimension is reduced from  $(m+n)$  to  $n$ .

In a discrete and finite  $\mathcal{S}^n$ , the convergence of ID-VI is guaranteed. Once  $V_{\pi}$  has converged to  $V_{\pi}^*$ , which is the optimal value function following policy  $\pi$ , we search for the predicted optimal state-to-reach  $\mathbf{x}_{t_{i+l}}^*$  by solving

$$\mathbf{x}_{t_{i+l}}^* = \arg \min_{\mathbf{x}_{t_{i+l}} \in \mathcal{S}^n} \left\{ \sum_{p=0}^{l-1} L(\mathbf{x}_{t_i+p}, \pi(\mathbf{x}_{t_i+p}, \mathbf{x}_{t_{i+l}}, \hat{\theta}_{t_i})) + V_{\pi}^*(\mathbf{x}_{t_{i+l}}) \right\}. \quad (9)$$

where  $l$  is the number of steps it takes to propagate from  $\mathbf{x}_{t_i}$  to  $\mathbf{x}_{t_{i+l}}$ , and  $l$  varies depending on  $\mathbf{x}_{t_{i+l}}$ . Because of the control constraints, some states in  $\mathcal{S}^n$  are not reachable from  $\mathbf{x}_{t_i}$  in one time step. Thus, the  $l$ -step lookahead strategy (Sutton, 1988) is used in (9) to evaluate those states as the state-to-reach. During the  $l$ -step lookahead, controls are constrained and the state constraints are introduced as penalty.

**Remark 3.1:** Because of the state constraints, e.g. collision avoidance and obstacle avoidance, the reachable set of states from  $\mathbf{x}_{t_i}$  following the ID policy will be different from the reachable set from  $\mathbf{x}_{t_{i+1}}$ . As a result, even if  $\mathbf{x}_{d,t_i} = \mathbf{x}_{d,t_{i+1}}$ , the predicted optimal state-to-reach at different time step will not necessarily be the same.

The ID-VI based control  $\mathbf{u}_{t_i}^*$  is calculated by substituting  $\mathbf{x}_{t_{i+l}}^*$  from (9) as the desired state-to-reach from  $\mathbf{x}_{t_i}$  into (7), i.e.

$$\mathbf{u}_{t_i}^* = \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+l}}^*, \hat{\theta}_{t_i}). \quad (10)$$

### 3.2 State space projection via VMC

To further reduce the dimension of the search space, the VMC rule (Srinivasan & Davey, 1995; Xu & Basset, 2012) is applied to confine the state search space to a VMC subspace  $\mathcal{S}_{v,t_i}^n \subseteq \mathcal{S}^n$  that is shaped by a reference state  $\mathbf{x}_{ref,t_i}$  and a desired state-to-reach  $\mathbf{x}_{d,t_i}$ . The VMC rule uniquely determines a state  $\mathbf{x}_{t_k}^v \in \mathcal{S}_{v,t_i}^n$  by a PCP  $v_{t_k}$  from the 1-dimensional PCP space  $\mathcal{V}^1$ . Following this procedure, searching for the optimal state-to-reach in the  $m$ -dimensional state space is transferred into searching for the optimal PCP-to-go in the 1-dimensional PCP space, which is then projected back into the state space to obtain the optimal state-to-reach.

Following the VMC rule, a  $\mathbf{x}_{t_k}^v \in \mathcal{S}_{v,t_i}^n$  can be determined by a PCP  $v_{t_k} \in \mathcal{V}^1$  as Xu and Basset (2012)

$$\begin{aligned} \mathbf{x}_{t_k}^v &= \mathbf{x}_{ref,t_i} + v_{t_k}(\mathbf{x}_{d,t_i} - \mathbf{x}_{ref,t_i}) \\ &\triangleq \mathbf{x}(v_{t_k}, \mathbf{x}_{d,t_i}, \mathbf{x}_{ref,t_i}), \quad v_{t_k} \in \mathcal{V}^1 \end{aligned} \quad (11)$$

where  $\mathbf{x}_{ref,t_i}$  corresponds to a PCP value of 0, and  $\mathbf{x}_{d,t_i}$  corresponds to a PCP value of 1. It is practical to let  $\mathcal{V}^1 = [0, v_{\max}]$ , where  $v_{\max}$  can be chosen based on desired system characteristics. For example,  $v_{\max} = 1$  corresponds to a desired state response without overshoot, while  $v_{\max} = 1.2$  allows for a 20% overshoot. Different from Xu and Basset (2012),  $\mathbf{x}_{t_i}$  does not necessarily belong to  $\mathcal{S}_{v,t_i}^n$ . The reference state  $\mathbf{x}_{ref,t_i}$  can be selected in such a way that  $\mathcal{S}_{v,t_i}^n \subseteq \mathcal{S}^n$ , i.e. every  $\mathbf{x}_{t_k}^v \in \mathcal{S}_{v,t_i}^n$  satisfies the state constraints. When the state constraints are too complicated to satisfy, desired waypoint states-to-go can be inserted to divide the original problem into sub-problems.

Equation (11) is a bijection, meaning for each  $v_{t_k} \in \mathcal{V}^1$ , there is a corresponding  $\mathbf{x}_{t_k}^v \in \mathcal{S}_{v,t_i}^n$ , and vice versa. The inverse projection from the VMC subspace  $\mathcal{S}_{v,t_i}^n$  to the PCP space  $\mathcal{V}^1$  is

$$\begin{aligned} v_{t_k} &= \langle \mathbf{x}_{t_k}^v - \mathbf{x}_{ref,t_i}, \mathbf{x}_{d,t_i} - \mathbf{x}_{ref,t_i} \rangle / \|\mathbf{x}_{d,t_i} - \mathbf{x}_{ref,t_i}\|_2 \\ &\triangleq v(\mathbf{x}_{t_k}^v, \mathbf{x}_{d,t_i}, \mathbf{x}_{ref,t_i}), \quad \mathbf{x}_{d,t_i} \neq \mathbf{x}_{ref,t_i}, \mathbf{x}_{t_k}^v \in \mathcal{S}_{v,t_i}^n. \end{aligned} \quad (12)$$

Substituting (11) into (7) with  $\mathbf{x}_{d,t_i}$  as the desired state-to-reach, the ID policy under the VMC rule becomes

$$\mathbf{u}_{\pi,t_k}^v = \pi(\mathbf{x}(v_{t_k}, \mathbf{x}_{d,t_i}, \mathbf{x}_{ref,t_i}), \mathbf{x}_{d,t_i}, \hat{\theta}_{t_i}). \quad (13)$$

The Bellman equation for ID-VMC-VI is obtained by substituting (11) and (13) into (8) as

$$\begin{aligned} V_{\pi,j}(\mathbf{x}(v_{t_k}, \mathbf{x}_{d,t_i}, \mathbf{x}_{ref,t_i})) \\ = L(\mathbf{x}_{t_k}^v, \mathbf{u}_{\pi,t_k}^v) + V_{\pi,j-1}(\mathbf{x}(v_{t_{k+1}}, \mathbf{x}_{d,t_i}, \mathbf{x}_{ref,t_i})), \\ \forall v_{t_k} \in \mathcal{V}^1, \end{aligned} \quad (14)$$

or simply

$$V_{\pi v,j}(v_{t_k}) = L(\mathbf{x}_{t_k}^v, \mathbf{u}_{\pi,t_k}^v) + V_{\pi v,j-1}(v_{t_{k+1}}), \quad \forall v_{t_k} \in \mathcal{V}^1 \quad (15)$$

where we use the additional subscript  $v$  to illustrate that the value function is obtained under the VMC rule.  $v_{t_{k+1}}$  is the PCP corresponding to  $F(\mathbf{x}_{t_k}^v, \hat{\theta}_{t_i}) + G(\mathbf{x}_{t_k}^v, \hat{\theta}_{t_i})\mathbf{u}_{\pi,t_k}^v$  and is calculated using (12).

ID-VMC-VI solves (15) for the optimal value function  $V_{\pi v}^*$  in the 1-dimensional PCP space. After this reduction step, the search space dimension is further reduced from  $n$  in ID-VI to 1.

For a discrete and finite  $\mathcal{V}^1$ , ID-VMC-VI is guaranteed to converge. Once  $V_{\pi v}^*$  is reached, we search  $\mathcal{V}^1$  for the predicted optimal PCP-to-go from  $\mathbf{x}_{t_i}$  as

$$v_{t_{i+l}}^* = \arg \min_{v_{t_{i+l}} \in \mathcal{V}^1} \left\{ \sum_{p=0}^{l-1} L(\mathbf{x}_{t_i+p}, \pi(\mathbf{x}_{t_i+p}, \mathbf{x}_{t_{i+l}}^v, \hat{\theta}_{t_i})) + V_{\pi v}^*(v_{t_{i+l}}) \right\} \quad (16)$$

where  $\mathbf{x}_{t_{i+l}}^v = \mathbf{x}(v_{t_{i+l}}, \mathbf{x}_{d,t_i}, \mathbf{x}_{ref,t_i})$ .  $l$  is the number of steps taken to reach  $\mathbf{x}_{t_{i+l}}^v$  from  $\mathbf{x}_{t_i}$ , and  $l$  is free. The  $l$ -step lookahead (Sutton, 1988) is warranted by not only control constraints, but



also the fact that  $\mathbf{x}_{t_i}$  may not belong to  $\mathcal{S}_{v,t_i}^n$ . Either way, the  $l$ -step lookahead is necessary for finding the optimal PCP-to-go over the 1-dimensional PCP space. Similar to ID-VI, the control and state constraints are considered in the  $l$ -step lookahead in (16).

**Remark 3.2:** Because of the state constraints, the optimal PCP-to-go changes as the system state propagates, even if the VMC subspace stays unchanged. Choosing  $v_{t_{i+l}}^*$  as the optimal PCP-to-go at  $t_i$  does not necessarily mean that the state trajectory will eventually reach  $\mathbf{x}(v_{t_{i+l}}^*, \mathbf{x}_{d,t_i}, \mathbf{x}_{ref,t_i})$  at  $t_{i+l}$ .

Substituting (16) into (11), the optimal  $l$ -step state-to-reach predicted by ID-VMC-VI is

$$\mathbf{x}_{t_{i+l}}^{v*} = \mathbf{x}(v_{t_{i+l}}^*, \mathbf{x}_{d,t_i}, \mathbf{x}_{ref,t_i}). \quad (17)$$

The ID-VMC-VI based control  $\mathbf{u}_{t_i}^{v*}$  is then calculated by substituting (17) as the desired state-to-reach from  $\mathbf{x}_{t_i}$  into (7), i.e.

$$\mathbf{u}_{t_i}^{v*} = \pi(\mathbf{x}_{t_i}, \mathbf{x}(v_{t_{i+l}}^*, \mathbf{x}_{d,t_i}, \mathbf{x}_{ref,t_i}), \hat{\boldsymbol{\theta}}_{t_i}). \quad (18)$$

### 3.3 Stability analysis

**Proposition 3.1:** At  $t_i$ , under the assumption that  $\hat{\boldsymbol{\theta}}_{t_i} = \boldsymbol{\theta}$ ,  $\mathbf{x}_{d,t_i}$  is the asymptotically stable point of the system (1) under control (18).

**Proof:** Let  $\mathbf{x}_{t_i} = \mathbf{x}_{d,t_i}$ . It is straight forward to get  $\mathbf{x}_{d,t_i} = F(\mathbf{x}_{d,t_i}, \hat{\boldsymbol{\theta}}_{t_i}) + G(\mathbf{x}_{d,t_i}, \hat{\boldsymbol{\theta}}_{t_i})\mathbf{u}_{t_i}^{v*}$ . Therefore,  $\mathbf{x}_{d,t_i}$  is an equilibrium point of the system (1) under control (18).

At  $t_i$ , the optimal value function  $V_{\pi v}^*$  is stationary. From (16), we define the value function of state  $\mathbf{x}_{t_i}$  under control (18) as

$$\begin{aligned} V_{\pi v}(\mathbf{x}_{t_i}) &\triangleq \sum_{p=0}^{l-1} L(\mathbf{x}_{t_{i+p}}, \pi(\mathbf{x}_{t_{i+p}}, \mathbf{x}_{t_{i+l}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}^*(v_{t_{i+l}}^*) \\ &= L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+l}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})) \\ &\quad + \sum_{p=1}^{l-1} L(\mathbf{x}_{t_{i+p}}, \pi(\mathbf{x}_{t_{i+p}}, \mathbf{x}_{t_{i+l}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}^*(v_{t_{i+l}}^*), \end{aligned} \quad (19)$$

in which

$$\begin{aligned} &\sum_{p=1}^{l-1} L(\mathbf{x}_{t_{i+p}}, \pi(\mathbf{x}_{t_{i+p}}, \mathbf{x}_{t_{i+l}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}^*(v_{t_{i+l}}^*) \\ &\stackrel{l_2=1}{=} \sum_{p=0}^{l_2-1} L(\mathbf{x}_{t_{i+1+p}}, \pi(\mathbf{x}_{t_{i+1+p}}, \mathbf{x}_{t_{i+1+l_2}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}^*(v_{t_{i+1+l_2}}^*) \\ &\geq \min_{v_{t_{i+1+l_3}} \in \mathcal{V}^1} \left\{ \sum_{p=0}^{l_3-1} L(\mathbf{x}_{t_{i+1+p}}, \pi(\mathbf{x}_{t_{i+1+p}}, \mathbf{x}_{t_{i+1+l_3}}^v, \hat{\boldsymbol{\theta}}_{t_i})) \right. \\ &\quad \left. + V_{\pi v}^*(v_{t_{i+1+l_3}}^*) \right\} \\ &= V_{\pi v}(\mathbf{x}_{t_{i+1}}). \end{aligned} \quad (20)$$

From the definition of the step cost function following (3), it is clear that  $V_{\pi v}(\mathbf{x}_{t_i}) \geq 0$ ,  $\forall \mathbf{x}_{t_i} \in \mathcal{S}^n$ , with  $V_{\pi v}(\mathbf{x}_{t_i}) = 0$  only when  $\mathbf{x}_{t_i} = \mathbf{x}_{d,t_i}$ .

Combining (19) and (20) yields

$$\begin{aligned} V_{\pi v}(\mathbf{x}_{t_{i+1}}) - V_{\pi v}(\mathbf{x}_{t_i}) &\leq -L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+l}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})) \leq 0, \\ \forall \mathbf{x}_{t_i} &\in \mathcal{S}^n. \end{aligned} \quad (21)$$

The second equality in (21) holds only when  $\mathbf{x}_{t_i} = \mathbf{x}_{d,t_i}$ . Therefore,  $V_{\pi v}(\mathbf{x}_{t_i})$  is a Lyapunov function for the system (1) under control (18) at  $t_i$ . Based on the discrete-time Lyapunov theorem (Luenberger, 1979) and the discrete-time LaSalle's invariance principle (Sundarapandian, 2003), the equilibrium point  $\mathbf{x}_{d,t_i}$  is asymptotically stable. ■

### 3.4 Optimality analysis

**Lemma 3.2:** When applied to solve C1,  $\mathbf{u}_{t_i}^{s*}$  is suboptimal to  $\mathbf{u}_{t_i}^*$ .

**Proof:** From (5), we have

$$\begin{aligned} V^*(\mathbf{x}_{t_i}) &= L(\mathbf{x}_{t_i}, \mathbf{u}_{t_i}^*) + V^*(F(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i}) + G(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})\mathbf{u}_{t_i}^*) \\ &\leq L(\mathbf{x}_{t_i}, \mathbf{u}_{t_i}^{s*}) + V^*(F(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i}) + G(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})\mathbf{u}_{t_i}^{s*}). \end{aligned} \quad (22)$$

Therefore,  $\mathbf{u}_{t_i}^{s*}$  is suboptimal to  $\mathbf{u}_{t_i}^*$ . ■

As mentioned earlier, the selection of reference point in the second dimension reduction step is mainly exploited for addressing the state constraints. For the purpose of optimality analysis,  $\mathbf{x}_{t_i}$  will be regarded as the reference point in this subsection.

**Lemma 3.3:** When applied to solve C1,  $\mathbf{u}_{t_i}^{v*}$  is suboptimal to  $\mathbf{u}_{t_i}^{s*}$ .

**Proof:** Choosing  $\mathbf{x}_{ref,t_i} = \mathbf{x}_{t_i}$ , then from (9), we have

$$\mathbf{x}_{t_{i+1}}^{s*} = \arg \min_{\mathbf{x}_{t_{i+1}} \in \mathcal{S}^n} \left\{ L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi}^*(\mathbf{x}_{t_{i+1}}) \right\}, \quad (23)$$

and from (16) and (17), we have

$$\begin{aligned} \mathbf{x}_{t_{i+1}}^{v*} &= \arg \min_{\mathbf{x}_{t_{i+1}}^{v*} \in \mathcal{S}_{v,t_i}^n} \left\{ L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^v, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}^*(\mathbf{x}_{t_{i+1}}^v) \right\} \\ &\geq \arg \min_{\mathbf{x}_{t_{i+1}}^{v*} \in \mathcal{S}_{v,t_i}^n} \left\{ L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^v, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi}^*(\mathbf{x}_{t_{i+1}}^v) \right\}. \end{aligned} \quad (24)$$

Because  $\mathcal{S}_{v,t_i}^n \subseteq \mathcal{S}^n$ , comparing (23) and (24) yields

$$\begin{aligned} &L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{s*}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi}^*(\mathbf{x}_{t_{i+1}}^{s*}) \\ &\leq L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}^*(\mathbf{x}_{t_{i+1}}^{v*}). \end{aligned} \quad (25)$$

From (10) and (18), we have  $\mathbf{u}_{t_i}^{s*} = \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{s*}, \hat{\boldsymbol{\theta}}_{t_i})$  and  $\mathbf{u}_{t_i}^{v*} = \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})$ . Therefore,  $\mathbf{u}_{t_i}^{v*}$  is suboptimal to  $\mathbf{u}_{t_i}^{s*}$ . ■

**Corollary 3.4:** For  $\forall \mathbf{x}_{t_{i+1}}^{v*} \in \mathcal{S}_{v,t_i}^n$ ,  $\pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})$  is suboptimal to  $\mathbf{u}_{t_i}^{v*}$  in solving C1.

**Proof:** The following inequality is a direct result from (24)

$$L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{v*}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}^*(\mathbf{x}_{t_{i+1}}^{v*}) \leq L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^v, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}^*(\mathbf{x}_{t_{i+1}}^v), \quad (26)$$

which proves the statement. ■

**Corollary 3.5:** If  $v \in [0, 1]$ , the optimality of using  $\mathbf{u}_{t_i}^{v*}$  in C1 is lower bounded by that of a first-order, exponential decay state tracking controller.

**Proof:** A first-order state trajectory in the VMC subspace can be formulated as Ogata (2010)

$$\mathbf{x}_{t_k} = \mathbf{x}_d - (\mathbf{x}_d - \mathbf{x}_i)e^{a(t_i - t_k)}, \quad k = i, \dots, N-1 \quad (27)$$

with  $1/a > 0$  being the time constant.

Substituting (27) into (12), the PCP-to-go can be derived as

$$v_{t_k} = 1 - e^{a(t_i - t_k)}, \quad k = i, \dots, N-1. \quad (28)$$

It is clear from (28) that  $0 < v_{t_{i+1}} < 1$ ,  $\forall a > 0$ . Based on Corollary 3.4, the tracking controller following this predefined PCP is suboptimal to  $\mathbf{u}_{t_i}^{v*}$ . ■

**Corollary 3.6:** If  $v$  can be larger than 1, the optimality of using  $\mathbf{u}_{t_i}^{v*}$  is lower bounded by that of a second-order state trajectory tracking controller.

**Proof:** A second-order state trajectory can be formulated as Ogata (2010)

$$\mathbf{x}_{t_k} = \mathbf{x}_d - (\mathbf{x}_d - \mathbf{x}_i) \frac{e^{\omega_d(t_i - t_k)\zeta/\sqrt{1-\zeta^2}}}{\sqrt{1-\zeta^2}} \sin \chi, \quad 0 < \zeta < 1, \quad k = i, \dots, N-1 \quad (29)$$

where  $\chi = \omega_d(t_k - t_i) + \tan^{-1}(\sqrt{1-\zeta^2}/\zeta)$ .  $\zeta$  is the damping ratio and  $\omega_d$  is the damped natural frequency.

The PCP-to-go following this second-order trajectory can be derived as

$$v_{t_k} = 1 - \frac{e^{\omega_d(t_i - t_k)\zeta/\sqrt{1-\zeta^2}}}{\sqrt{1-\zeta^2}} \cdot [\sin \omega_d(t_k - t_i) \cot \chi + \cos \chi], \quad 0 < \zeta < 1, \quad k = i, \dots, N-1. \quad (30)$$

Choosing  $k = i+1$  in (30) gives the PCP-to-go  $v_{t_{i+1}}$ . Following Corollary 3.4, a controller tracking this predefined PCP is suboptimal to  $\mathbf{u}_{t_i}^{v*}$ . ■

**Remark 3.3:** The performance advantage of the proposed ID-VMC-VI control scheme over typical trajectory tracking controllers lies in the optimisation over PCPs. Furthermore, in ID-VMC-VI, we have the flexibility to choose the reference point  $\mathbf{x}_{ref, t_i}$  to address state constraints.

### 3.5 Dimension analysis

Following the two-step search space dimension reduction strategy, the dimension of search space for the optimal control is reduced from  $m+n$  in basic VI to  $n$  in ID-VI, and eventually reduced to 1 in ID-VMC-VI.

Assume that there are  $D$  discretised nodes per dimension in the search space. Because the size of a tabular value function is proportional to the number of discretised nodes in the search space, the memory consumption of ID-VMC-VI will be  $D^{n-1}$  and  $D^{m+n-1}$  times less than those of ID-VI and VI, respectively.

Even though the computational cost of a VI algorithm depends on many factors besides the size of the search space and no rigid relationship can be drawn between the two, ID-VMC-VI still dramatically reduces the computational cost compared with ID-VI and the basic VI, as will be shown in the simulation results.

## 4. ID policy improvement

### 4.1 Uncertain parameter adaptation via gradient descent

**E1:** The problem of estimating the uncertain parameter  $\boldsymbol{\theta}$  to improve the ID policy, as stated in Section 2, is solved by the gradient descent method (Curry, 1944)

$$\hat{\boldsymbol{\theta}}_{t_{i+1}} = \hat{\boldsymbol{\theta}}_{t_i} - \alpha \frac{de_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} \quad (31)$$

where  $\alpha$  is the step size matrix.  $de_{t_{i+1}}/d\hat{\boldsymbol{\theta}}_{t_i}$  is derived as

$$\frac{de_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} = \frac{d\tilde{\mathbf{x}}_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} \frac{de_{t_{i+1}}}{d\tilde{\mathbf{x}}_{t_{i+1}}} = \frac{d\tilde{\mathbf{x}}_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} \tilde{\mathbf{x}}_{t_{i+1}} \quad (32)$$

where

$$\frac{d\tilde{\mathbf{x}}_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} = \frac{d\mathbf{x}_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} - \frac{d\hat{\mathbf{x}}_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}}. \quad (33)$$

Given  $\mathbf{x}_{t_i}$ ,  $\mathbf{x}_{t_{i+1}} = F(\mathbf{x}_{t_i}, \boldsymbol{\theta}) + G(\mathbf{x}_{t_i}, \boldsymbol{\theta})\mathbf{u}_{t_i}$ , where  $\mathbf{u}_{t_i} = G^{-1}(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})[\hat{\mathbf{x}}_{t_{i+1}} - F(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})]$ . The derivative of  $\mathbf{x}_{t_{i+1}}$  with respect to  $\hat{\boldsymbol{\theta}}_{t_i}$  via the chain rule is

$$\begin{aligned} \frac{d\mathbf{x}_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} &= \frac{d\mathbf{u}_{t_i}}{d\hat{\boldsymbol{\theta}}_{t_i}} \frac{d\mathbf{x}_{t_{i+1}}}{d\mathbf{u}_{t_i}} \\ &= \left( \frac{\partial \mathbf{u}_{t_i}}{\partial \hat{\boldsymbol{\theta}}_{t_i}} + \frac{d\hat{\mathbf{x}}_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} \frac{\partial \mathbf{u}_{t_i}}{\partial \hat{\mathbf{x}}_{t_{i+1}}} \right) \frac{d\mathbf{x}_{t_{i+1}}}{d\mathbf{u}_{t_i}} \end{aligned} \quad (34)$$

where

$$\begin{aligned} \frac{\partial \mathbf{u}_{t_i}}{\partial \hat{\boldsymbol{\theta}}_{t_i}} &= \frac{dG^{-1}(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})}{d\hat{\boldsymbol{\theta}}_{t_i}} [\hat{\mathbf{x}}_{t_{i+1}} - F(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})] \\ &\quad - \left[ \frac{dF(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})}{d\hat{\boldsymbol{\theta}}_{t_i}} \right]^T [G^{-1}(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})]^T. \end{aligned} \quad (35)$$

On the other hand,  $\hat{\mathbf{x}}_{t_{i+1}}$  satisfies  $\hat{\mathbf{x}}_{t_{i+1}} = F(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i}) + G(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})\mathbf{u}_{t_i}$ , into which substituting  $\mathbf{u}_{t_i} = G^{-1}(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})[\hat{\mathbf{x}}_{t_{i+1}} - F(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})]$

yields  $\hat{\mathbf{x}}_{t_{i+1}} = \hat{\mathbf{x}}_{t_i+1}$ . Therefore,  $\hat{\mathbf{x}}_{t_{i+1}}$  is invariant with respect to  $\hat{\boldsymbol{\theta}}_{t_i}$ , i.e.  $d\hat{\mathbf{x}}_{t_{i+1}}/d\hat{\boldsymbol{\theta}}_{t_i} = \mathbf{0}$ . Hence, from (33),

$$d\tilde{\mathbf{x}}_{t_{i+1}}/d\hat{\boldsymbol{\theta}}_{t_i} = d\mathbf{x}_{t_{i+1}}/d\hat{\boldsymbol{\theta}}_{t_i}, \quad (36)$$

and from (34),

$$\frac{d\mathbf{x}_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} = \frac{\partial \mathbf{u}_{t_i}}{\partial \hat{\boldsymbol{\theta}}_{t_i}} \frac{d\mathbf{x}_{t_{i+1}}}{d\mathbf{u}_{t_i}}. \quad (37)$$

Combining (37) and (36), we have

$$\frac{d\tilde{\mathbf{x}}_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} = \frac{\partial \mathbf{u}_{t_i}}{\partial \hat{\boldsymbol{\theta}}_{t_i}} \frac{d\mathbf{x}_{t_{i+1}}}{d\mathbf{u}_{t_i}}. \quad (38)$$

Substituting (38) into (32), we have

$$\frac{de_{t_{i+1}}}{d\hat{\boldsymbol{\theta}}_{t_i}} = \frac{\partial \mathbf{u}_{t_i}}{\partial \hat{\boldsymbol{\theta}}_{t_i}} \frac{d\mathbf{x}_{t_{i+1}}}{d\mathbf{u}_{t_i}} \tilde{\mathbf{x}}_{t_{i+1}}, \quad (39)$$

which, when substituted into (31), yields

$$\hat{\boldsymbol{\theta}}_{t_{i+1}} = \hat{\boldsymbol{\theta}}_{t_i} - \alpha \frac{\partial \mathbf{u}_{t_i}}{\partial \hat{\boldsymbol{\theta}}_{t_i}} \frac{d\mathbf{x}_{t_{i+1}}}{d\mathbf{u}_{t_i}} \tilde{\mathbf{x}}_{t_{i+1}}. \quad (40)$$

## 4.2 Convergence analysis

**Proposition 4.1:** *Following the gradient method (40), the adaptation parameter  $\hat{\boldsymbol{\theta}}$  converges to the actual value  $\boldsymbol{\theta}$  if the system is in a transient stage.*

**Proof:** According to Curry (1944), the gradient descent method stops only when a stationary point is reached, i.e. when

$$de_{t_{i+1}}/d\hat{\boldsymbol{\theta}}_{t_i} = \mathbf{0}. \quad (41)$$

Given the fact that the system state is affected by the uncertain parameter,  $d\tilde{\mathbf{x}}_{t_{i+1}}/d\hat{\boldsymbol{\theta}}_{t_i} \neq \mathbf{0}$ . Therefore, we know from (32) that the stopping condition (41) holds true only if  $\tilde{\mathbf{x}}_{t_{i+1}} = \mathbf{0}$ . For a system in transient stages, this condition can be elaborated as

$$\begin{aligned} \tilde{\mathbf{x}}_{t_{i+1}} &= [F(\mathbf{x}_{t_i}, \boldsymbol{\theta}) - F(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})] \\ &\quad + [G(\mathbf{x}_{t_i}, \boldsymbol{\theta})\mathbf{u}_{t_i} - G(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i})\mathbf{u}_{t_i}] = \mathbf{0}, \\ \forall \mathbf{x}_{t_i} \in \mathcal{S}^n, \end{aligned} \quad (42)$$

For a general system, we assume there are an infinite number of states in the state space, therefore, Equation (42) holds true only if  $\hat{\boldsymbol{\theta}}_{t_i} = \boldsymbol{\theta}$ . ■

## 4.3 Stability condition

To discuss the influence of parameter adaptation on the stability of the closed-loop system,  $V(\cdot, \hat{\boldsymbol{\theta}})$  is used to denote a value function attained under the estimate  $\hat{\boldsymbol{\theta}}$ . Definition of  $V_{\pi v}(\mathbf{x}, \hat{\boldsymbol{\theta}})$  follows (19).

**Proposition 4.2:** *Assuming that there exists a vector  $\boldsymbol{\kappa} \in \mathbb{R}^r$ , such that  $\max_{\mathbf{x}_{t_k} \in \mathcal{S}^n} [V_{\pi v}(\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_{i+1}}) - V_{\pi v}(\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_i})] \leq \boldsymbol{\kappa}^T (\hat{\boldsymbol{\theta}}_{t_{i+1}} - \hat{\boldsymbol{\theta}}_{t_i})$ ,  $k \geq i$ , then  $(\mathbf{x}_{d,t_i}, \boldsymbol{\theta})$  is the asymptotically stable*

*point of the system (1) under control (18) and parameter adaptation (40) if the parameter update step size  $\alpha$  satisfies*

$$\boldsymbol{\kappa}^T \alpha \frac{\partial \mathbf{u}_{t_i}}{\partial \hat{\boldsymbol{\theta}}_{t_i}} \frac{d\mathbf{x}_{t_{i+1}}}{d\mathbf{u}_{t_i}} \tilde{\mathbf{x}}_{t_{i+1}} + L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{u*}, \hat{\boldsymbol{\theta}}_{t_i})) \geq 0 \quad (43)$$

*with the equality holds only when  $\mathbf{x}_{t_i} = \mathbf{x}_{d,t_i}$ .*

**Proof:** Proof of Proposition 4.2 follows a similar idea of Proposition 3.1. It is clear from Propositions 3.1 and 4.1 that  $(\mathbf{x}_{d,t_i}, \boldsymbol{\theta})$  is an equilibrium point of the overall system.

From the definition of the step cost function following (3),  $V_{\pi v}(\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_h}) \geq 0$ ,  $\forall (\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_h}) \in \{(\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_h}) | \mathbf{x}_{t_k} \in \mathcal{S}^n, \hat{\boldsymbol{\theta}}_{t_h} \in \mathcal{P}^r\}$ ,  $k \geq i, h \geq i$ .  $V_{\pi v}(\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_h}) = 0$  only when  $\mathbf{x}_{t_k} = \mathbf{x}_{d,t_i}$ .

From (19), we have

$$\begin{aligned} V_{\pi v}(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i}) &= L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{u*}, \hat{\boldsymbol{\theta}}_{t_i})) \\ &\quad + \sum_{p=1}^{l-1} L(\mathbf{x}_{t_{i+p}}, \pi(\mathbf{x}_{t_{i+p}}, \mathbf{x}_{t_{i+l}}^{u*}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}^*(\mathbf{u}_{t_{i+l}}^*) \\ &= L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{u*}, \hat{\boldsymbol{\theta}}_{t_i})) + V_{\pi v}(\mathbf{x}_{t_{i+1}}, \hat{\boldsymbol{\theta}}_{t_i}). \end{aligned} \quad (44)$$

From (44) and the assumption of Proposition 4.2, we have

$$\begin{aligned} &V_{\pi v}(\mathbf{x}_{t_{i+1}}, \hat{\boldsymbol{\theta}}_{t_{i+1}}) - V_{\pi v}(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i}) \\ &= V_{\pi v}(\mathbf{x}_{t_{i+1}}, \hat{\boldsymbol{\theta}}_{t_{i+1}}) - V_{\pi v}(\mathbf{x}_{t_{i+1}}, \hat{\boldsymbol{\theta}}_{t_i}) \\ &\quad - L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{u*}, \hat{\boldsymbol{\theta}}_{t_i})) \\ &\leq \boldsymbol{\kappa}^T (\hat{\boldsymbol{\theta}}_{t_{i+1}} - \hat{\boldsymbol{\theta}}_{t_i}) - L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{u*}, \hat{\boldsymbol{\theta}}_{t_i})). \end{aligned} \quad (45)$$

Substituting (40) into (45), we have

$$\begin{aligned} &V_{\pi v}(\mathbf{x}_{t_{i+1}}, \hat{\boldsymbol{\theta}}_{t_{i+1}}) - V_{\pi v}(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i}) \\ &\leq -\boldsymbol{\kappa}^T \alpha \frac{\partial \mathbf{u}_{t_i}}{\partial \hat{\boldsymbol{\theta}}_{t_i}} \frac{d\mathbf{x}_{t_{i+1}}}{d\mathbf{u}_{t_i}} \tilde{\mathbf{x}}_{t_{i+1}} - L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{u*}, \hat{\boldsymbol{\theta}}_{t_i})), \end{aligned} \quad (46)$$

from which it is clear that (43) is the condition for  $V_{\pi v}(\mathbf{x}_{t_{i+1}}, \hat{\boldsymbol{\theta}}_{t_{i+1}}) - V_{\pi v}(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i}) \leq 0$ . When  $\mathbf{x}_{t_i} = \mathbf{x}_{d,t_i}$ ,  $L(\mathbf{x}_{t_i}, \pi(\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}^{u*}, \hat{\boldsymbol{\theta}}_{t_i})) = 0$ . Thus, based on Proposition 4.1,  $V_{\pi v}(\mathbf{x}_{t_{i+1}}, \hat{\boldsymbol{\theta}}_{t_{i+1}}) - V_{\pi v}(\mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i}) = 0$  only if  $\hat{\boldsymbol{\theta}}_{t_i} = \boldsymbol{\theta}$ . Therefore,  $V_{\pi v}(\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_h})$ ,  $\forall (\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_h}) \in \{(\mathbf{x}_{t_k}, \hat{\boldsymbol{\theta}}_{t_h}) | \mathbf{x}_{t_k} \in \mathcal{S}^n, \hat{\boldsymbol{\theta}}_{t_h} \in \mathcal{P}^r\}$ ,  $k \geq i, h \geq i$ , is a Lyapunov function for the system (1) under control (18) and parameter adaptation (40), and the equilibrium point  $(\mathbf{x}_{d,t_i}, \boldsymbol{\theta})$  is asymptotically stable (Luenberger, 1979; Sundarapandian, 2003). ■

## 5. Algorithm

The ID-VMC-VI algorithm is summarised in Algorithm 1. All related equations have been discussed before. It is worth mentioning that the algorithm needs to run at every time step  $t_i$  to take time-varying properties, such as state constraints, into consideration, as well as to update the estimate of uncertain parameter.

There are several guidelines to help efficiently implement the algorithm. (1) A small step size is preferred to achieve better

**Algorithm 1** ID-VMC-VI Algorithm at  $t_i$ 


---

```

1: Given  $\mathbf{x}_{d,t_i}, \mathbf{x}_{t_i}, \hat{\boldsymbol{\theta}}_{t_i}, \boldsymbol{\alpha}, \mathcal{V}^1$ 
2: Set  $V_{\pi v,0}(v) = 0, \forall v \in \mathcal{V}^1, j = 1, \text{flag} = 0$ 
3: while  $\text{flag} = 0$  do
4:   for  $\forall v_{t_k} \in \mathcal{V}^1 \setminus \{v_d\}$  do
5:     Calculate  $\mathbf{x}_{t_k}^v = \mathbf{x}(v_{t_k}, \mathbf{x}_{d,t_i}, \mathbf{x}_{\text{ref},t_i})$ 
6:     Calculate  $\mathbf{u}_{\pi,t_k}^v = \pi(\mathbf{x}_{t_k}^v, \mathbf{x}_{d,t_i}, \hat{\boldsymbol{\theta}}_{t_i})$ 
7:     Propagate  $\mathbf{x}_{t_{k+1}}^v = F(\mathbf{x}_{t_k}^v, \hat{\boldsymbol{\theta}}_{t_i}) + G(\mathbf{x}_{t_k}^v, \hat{\boldsymbol{\theta}}_{t_i})\mathbf{u}_{\pi,t_k}^v$ 
8:     Calculate  $v_{t_{k+1}} = v(\mathbf{x}_{t_{k+1}}^v, \mathbf{x}_{d,t_i}, \mathbf{x}_{\text{ref},t_i})$ 
9:      $V_{\pi v,j}(v_{t_k}) \leftarrow L(\mathbf{x}_{t_k}^v, \mathbf{u}_{\pi,t_k}^v) + V_{\pi v,j-1}(v_{t_{k+1}})$ 
10:   end for
11:   if  $V_{\pi v,j}(v) = V_{\pi v,j-1}(v), \forall v \in \mathcal{V}^1$  then
12:     Set  $\text{flag} = 1$ 
13:   else
14:      $j \leftarrow j + 1$ 
15:   end if
16: end while
17: Calculate  $v_{t_{i+1}}^*$  following (16)
18: Calculate  $\mathbf{x}_{t_{i+1}}^{v*} = \mathbf{x}(v_{t_{i+1}}^*, \mathbf{x}_{d,t_i}, \mathbf{x}_{\text{ref},t_i})$ 
19: Calculate  $\mathbf{u}_{t_i}^{v*} = \pi(\mathbf{x}_{t_i}, \mathbf{x}(v_{t_{i+1}}^*, \mathbf{x}_{d,t_i}, \mathbf{x}_{\text{ref},t_i}), \hat{\boldsymbol{\theta}}_{t_i})$ 
20: Observe state feedback  $\mathbf{x}_{t_{i+1}}$ 
21: Calculate  $\tilde{\mathbf{x}}_{t_{i+1}} = \mathbf{x}_{t_{i+1}} - \mathbf{x}_{t_{i+1}}^{v*}$ 
22: Update  $\hat{\boldsymbol{\theta}}_{t_{i+1}} \leftarrow \hat{\boldsymbol{\theta}}_{t_i} - \boldsymbol{\alpha} \nabla_{\hat{\boldsymbol{\theta}}_{t_i}} e_{t_{i+1}}$ 

```

---

performance near the border of state and control constraints. (2) The solution optimality depends on the selection of reference point and prey motion (i.e. desired state  $\mathbf{x}_{d,t_i}$ ). In general, the reference point and desired state  $\mathbf{x}_{d,t_i}$  are chosen so that the VMC subspace does not violate the state constraints. However, in case of complicated state constraints, an intermediate desired state-to-reach can be adopted to break down the problem into sub-problems. (3) Control constraints should be accounted for in the state propagation in Line 9 of the above algorithm. (4) Both state and control constraints should be accounted for in the  $l$ -step lookahead in Line 14 of the algorithm.

**Remark 5.1:** We reckon that the ID-VMC-VI algorithm can be implemented in real-time if the CPU time to execute the algorithm is less than the control time step size.

## 6. Simulation validation and discussion

Two examples are simulated to show the computational cost and memory usage advantages of the ID-VMC-VI method. Example 1 solves a regulating problem with control constraints. Example 2 solves an optimal trajectory planning problem with both state and control constraints. While the first example is adopted and refined from Li and Xu (2020), the improved algorithm and the new ID policy improvement method from this paper are applied. Furthermore, the simulation parameters and settings are tuned to achieve better performance. Since the estimates of uncertain parameters are updated at every time step, the VI algorithm needs to be implemented in real time. Theoretically, the classical DP based VI method (VI) will generate the optimal solutions, however, its high computational cost and a huge amount of memory requirement prohibit it from real-time

application. Therefore, we show the comparison of VI, ID-VI, and ID-VMC-VI methods only in Example 1; while in Example 2, ID-VI and ID-VMC-VI are compared. The implementation of the ID-VI method has two differences from the ID-VMC-VI method. First, during the value iteration, the VMC rule is not followed. Second, (9) is followed instead of (16) to calculate the predicted optimal state-to-reach.

### 6.1 Example 1: angular velocity control

Simulation 1 is conducted on a laptop computer with a 3.7 GHz CPU and a 64 GB RAM. This example is adopted and refined from Li and Xu (2020). A UAV attitude dynamics is governed by Mahony et al. (2012)

$$\boldsymbol{\omega}_{t_{i+1}} = \boldsymbol{\omega}_{t_i} - T\mathbf{I}^{-1}(\boldsymbol{\omega}_{t_i} \times \mathbf{I}\boldsymbol{\omega}_{t_i}) + T\mathbf{I}^{-1}\mathbf{u}_{t_i} \quad (47)$$

where  $\boldsymbol{\omega}_{t_i} \in \mathcal{S}^3$  is the UAV body frame angular velocity.  $\mathbf{u}_{t_i} \in \mathcal{A}^3$  is the control torque imposed on the centre of mass of the UAV. Subscripts  $x, y$  and  $z$  are used to denote the components of  $\boldsymbol{\omega}_{t_i}$  and  $\mathbf{u}_{t_i}$  along the directions of roll, pitch, and yaw.  $\mathbf{I}$  is the moment of inertia matrix. It is invertible, but not perfectly known.  $T$  is the time step size. We assume the initial guess of  $\mathbf{I}$  to be  $\hat{\mathbf{I}} = \text{diag}\{0.07, 0.07, 0.1\} \text{ kg}\cdot\text{m}^2$  and the real values to be

$$\mathbf{I} = \begin{bmatrix} 0.09 & 0 & 0.01 \\ 0 & 0.05 & 0 \\ 0.01 & 0 & 0.13 \end{bmatrix} \text{ kg}\cdot\text{m}^2. \quad (48)$$

The proposed ID-VMC-VI algorithm is tasked to follow the step angular velocity command. The step cost function is

$$L(\boldsymbol{\omega}_{t_i}, \mathbf{u}_{t_i}) = (\boldsymbol{\omega}_{d,t_i} - \boldsymbol{\omega}_{t_i})^T \mathbf{Q}(\boldsymbol{\omega}_{d,t_i} - \boldsymbol{\omega}_{t_i}) + \mathbf{u}_{t_i}^T \mathbf{R} \mathbf{u}_{t_i} \quad (49)$$

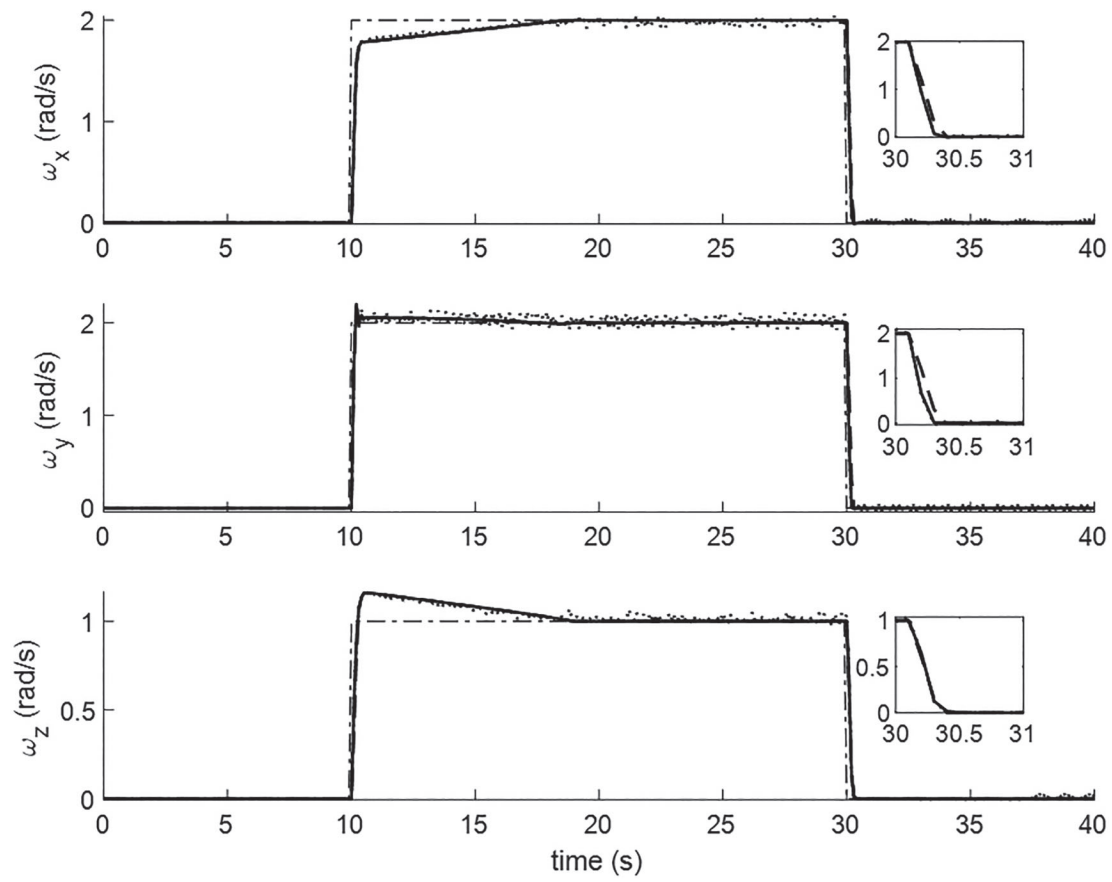
where the weighting matrices  $\mathbf{Q}$  and  $\mathbf{R}$  are picked to be  $\text{diag}\{1, 1, 1\}$  and  $\text{diag}\{0.01, 0.01, 0.01\}$ . To test the ability of the proposed method in handling control constraints, the control limits are assumed to be  $\pm 0.8 \text{ Nm}$  in all three directions.

The time step size is set to be 0.1 s in the simulation. For the purpose of comparison, the 6-dimensional state and action spaces of VI, 3-dimensional state space of ID-VI, and the 1-dimensional PCP space of ID-VMC-VI are discretised into  $41^6$  nodes,  $41^3$  nodes, and 41 nodes, respectively.

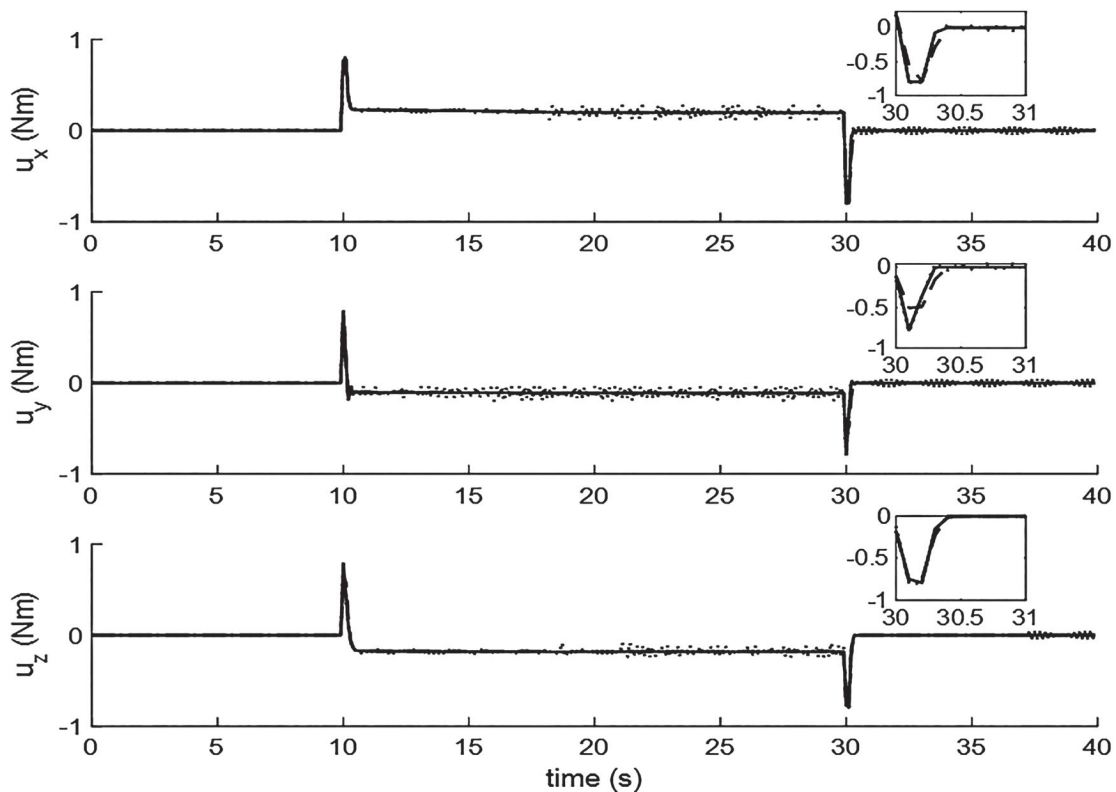
The step signals last 20 s in all three directions. In the roll and pitch directions, the signal magnitudes are 2 rad/s. In the yaw direction, the magnitude is 1 rad/s. The simulation results show that all three methods can stabilise the angular velocities (shown in Figure 1), without violating the control constraints (shown in Figure 2). The benefit of ID-VI and VI being able to search the state space and state/action spaces yields a performance boost over ID-VMC-VI at the rising phase of the step response (shown in the zoomed plots in Figure 1). As the errors in  $\hat{\mathbf{I}}$  converge to zeros (shown in Figure 3), the steady state errors in the step responses gradually decrease (shown in Figure 1), indicating that the control performance is being improved.

Table 1 shows the CPU time, memory usage, and performance indices for all three algorithms. On average, ID-VMC-VI is 654 times and 38,080 times faster than ID-VI and VI, respectively. It also uses 1607 times less memory with respect to ID-VI,

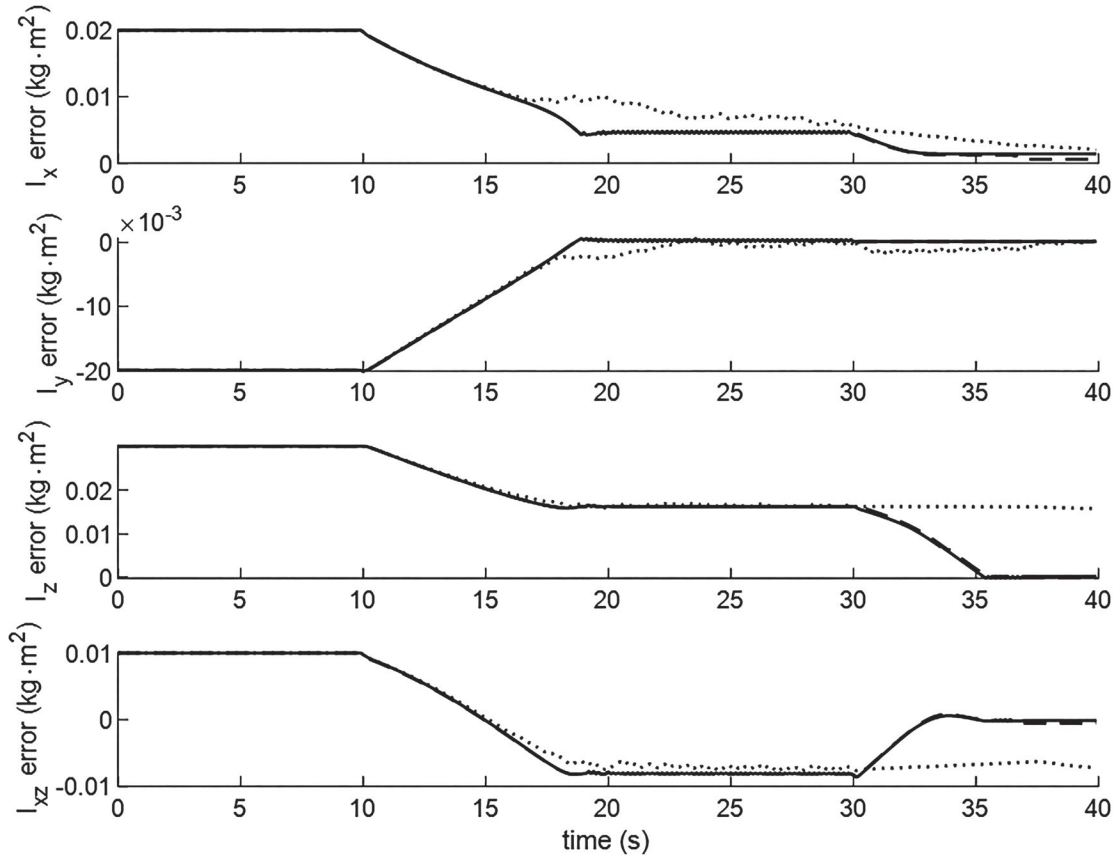




**Figure 1.** Comparison of the angular velocity step responses. (Dash-dotted lines: step signals; solid lines: ID-VI resulted step responses; dashed lines: ID-VMC-VI resulted step responses; and dotted lines: VI resulted step responses.)



**Figure 2.** Comparison of the control commands in response to the step signals. (Solid lines: ID-VI control commands; dashed lines: ID-VMC-VI control commands; and dotted lines: VI control commands.)



**Figure 3.** Comparison of the errors in estimation. (Solid lines: ID-VI; dashed lines: ID-VMC-VI; and dotted lines: VI)

**Table 1.** CPU time, memory usage and performance indices in Example 1.

Method	VI	ID-VI	ID-VMC-VI
Average CPU Time per simulation step (0.1 s)	23.8 s	0.409 s	0.625 ms
Memory usage	35.4 GB	3.15 MB	1.96 KB
Performance index	2.516	2.447	2.631

which is consistent with the analysis in Section 3.5. The memory usage of VI is huge, and is the main reason for a computer with a 64 GB RAM to be involved. This type of RAM requirement makes the VI algorithm practically infeasible for most mobile robotic platforms. In terms of performance index, the lesser the better. The performance index of ID-VMC-VI is 7.5% higher than that of ID-VI. Theoretically, VI should achieve the minimum performance index. However, practically more optimisable parameters may lead to challenges in convergence and the performance index might be higher. As we argued in Remark 5.1, it can be implemented in real time. One way to improve the performance of both ID-VI and ID-VMC-VI is to increase the number of discretised nodes in the search spaces. In contrast to the memory usage of ID-VI that grows exponentially with respect to the number of nodes in the search space, the memory usage of ID-VMC-VI is proportional to that number, which means that ID-VMC-VI is much more scalable. Thus, it can be implemented on platforms with limited memory resources and can be easily scaled up to provide a better performance.

## 6.2 Example 2: skid-steer vehicle minimum time-to-reach problem

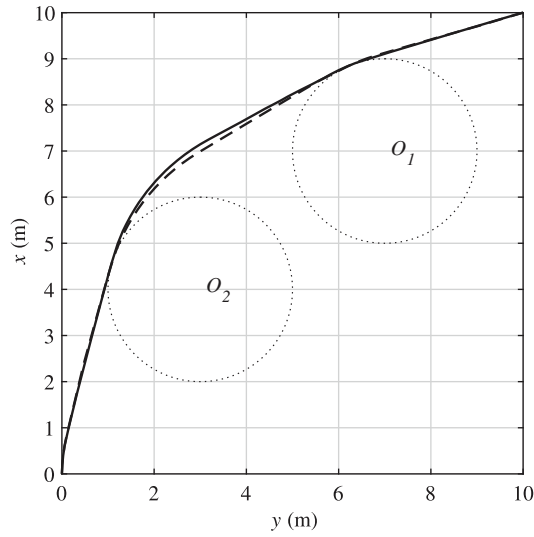
Simulation 2 is conducted on a laptop computer (CPU: 2.6 GHz; memory: 16 GB). The point mass model of the vehicle is modified from Laumond et al. (1998) as

$$\begin{aligned} x_{t_{i+1}} &= x_{t_i} + TV_{t_i} \cos \theta_{t_i}; & y_{t_{i+1}} &= y_{t_i} + TV_{t_i} \sin \theta_{t_i}; \\ \theta_{t_{i+1}} &= \theta_{t_i} + T(\omega_{t_i} + b_c) \end{aligned} \quad (50)$$

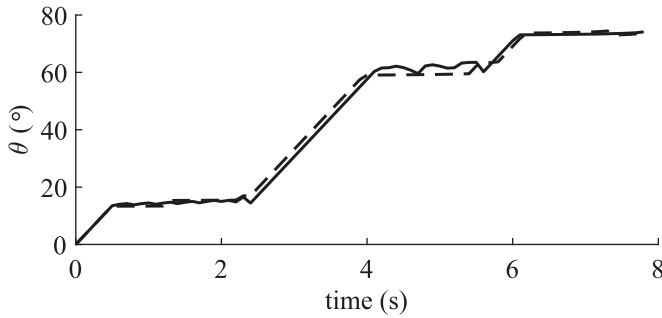
where  $(x_{t_i}, y_{t_i})$  is the vehicle location in the  $x$ - $y$  plane of an inertial north-east-down frame, of which the origin is on the ground surface.  $V_{t_i}$  is the velocity, while  $\theta_{t_i}$  and  $\omega_{t_i}$  are the heading angle and angular velocity around the  $z$  axis.  $T = 0.1$  s is the time step size. With the vehicle initially facing north,  $\theta_0 = 0^\circ$ . A constant control input bias  $b_c = -3^\circ/\text{s}$  is considered and estimated following the gradient based adaptation rule discussed in Section 4. The following constraints are used:  $0 \leq V \leq 2$  m/s and  $|\omega| \leq 30^\circ/\text{s}$ .

The proposed ID-VMC-VI algorithm is used to control the vehicle from the starting point (0, 0) to the target point (10, 10), while avoiding two circular obstacle regions  $O_1$  and  $O_2$ , which are respectively centred at (7, 7) and (4, 3) with a radius of 2 m. The step cost function is  $L = T$ , where  $T > 0$  is the constant time step size.

For the purpose of comparison, the 2-dimensional state space of ID-VI and the 1-dimensional PCP space of ID-VMC-VI are discretised into  $100 \times 100$  nodes and 100 nodes. The



**Figure 4.** Comparison of the optimised paths. (Dotted line: obstacle; solid line: ID-VI optimised path; dashed line: ID-VMC-VI optimised path.)



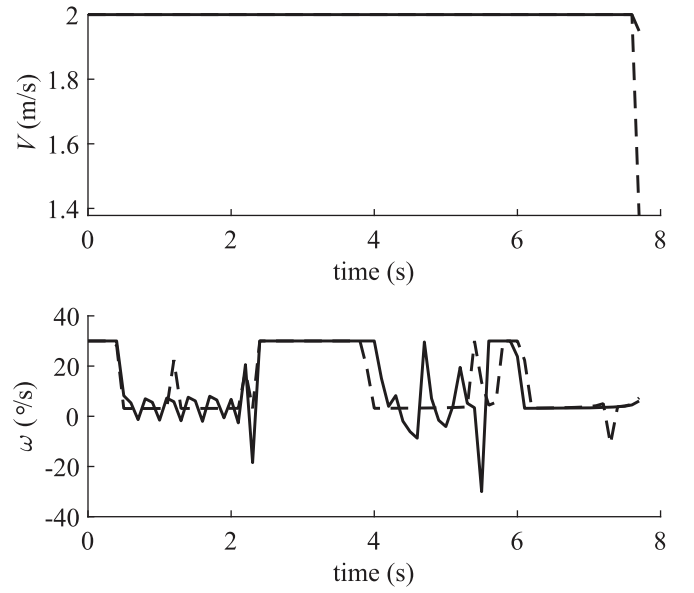
**Figure 5.** Comparison of the heading angles of the vehicle. (Solid line: ID-VI resulted heading angle; dashed line: ID-VMC-VI resulted heading angle.)

point where the upper tangent line of  $O_1$  through the target point intersects  $y = 0$  is used as the reference point.

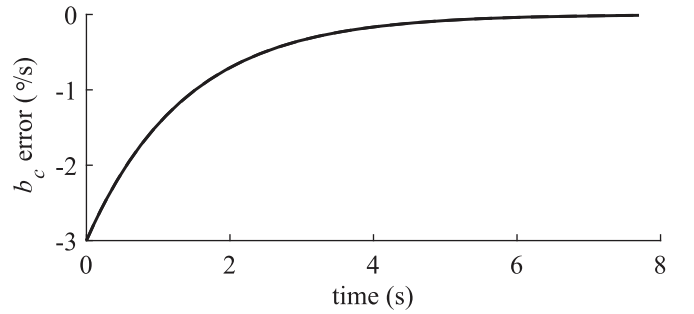
As shown in Figures 4 and 5, the ID-VI and ID-VMC-VI methods generate very similar paths that can avoid the obstacles and reach the target location. The ID-VI method introduces oscillations in the heading angle control due to the discretization of 2-dimensional state space. On the other hand, the ID-VMC-VI method generates much smoother heading angle control commands as the benefit of searching the 1-dimensional PCP space.

On the other hand, as shown in Figure 6, both methods generate control signals that satisfy the control constraints. Figure 7 shows that the adopted gradient based adaptation rule can efficiently reduce the estimation error in the input bias following both methods.

Table 2 shows the comparison of CPU time, memory usage and performance between ID-VI and ID-VMC-VI. On average, the ID-VMC-VI algorithm is 194 times faster to compute and consumes 100 times less memory than ID-VI, while delivering about the same performance index in this example.



**Figure 6.** Comparison of the control variables. (Solid lines: ID-VI control commands; dashed lines: ID-VMC-VI control commands.)



**Figure 7.** Comparison of the errors in the input bias estimation. (Solid lines: ID-VI resulted estimation error changes; dashed lines: ID-VMC-VI resulted estimation error changes. The solid line and the dashed line are almost identical.)

**Table 2.** CPU time, memory usage and performance indices in Example 2.

Method	ID-VI	ID-VMC-VI
Average CPU Time (s)	5.652	0.029
per simulation step (0.1 s)		
Memory usage (kilobyte)	234.375	2.344
Performance index	7.8	7.8

## 7. Conclusion

In this study, a dynamic programming search space dimension reduction strategy is proposed for a class of discrete-time dynamic systems. Two reduction steps are involved in the proposed value iteration algorithm: (1) eliminating the action search space using inverse dynamics; and (2) projecting the state search space into a 1-dimensional parameter space by the virtual motion camouflage rule. A gradient based uncertain parameter adaptation rule is developed to improve the inverse dynamics policy at each time step. Compared with typical dynamic programming based optimal controls, the proposed algorithm can dramatically reduce memory usage and computational time, and is thus feasible for real-time applications. Stability and optimality of the algorithm are analysed and its performance lower

bounds are given. Two simulation examples are used to illustrate the aforementioned salient features of the proposed algorithm.

## Disclosure statement

No potential conflict of interest was reported by the author(s).

## Funding

This work was supported by the National Science Foundation [grant number 1924622].

## References

- Barto, A. G., Sutton, R. S., & Anderson, C. W. (1983, September–October). Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-13(5), 834–846. <https://doi.org/10.1109/TSMC.1983.6313077>
- Bellman, R. (1954). Dynamic programming and a new formalism in the calculus of variations. *Proceedings of the National Academy of Sciences of the United States of America*, 40(4), 231–235. <https://doi.org/10.1073/pnas.40.4.231>
- Bellman, R. (1961). *Adaptive control processes: A guided tour*. Princeton University Press.
- Bertsekas, D. P. (2013). Rollout algorithms for discrete optimization: A survey. In P. M. Pardalos, D. Du, & R. L. Graham (Eds.), *Handbook of combinatorial optimization* (pp. 2989–3013). Springer.
- Bertsekas, D. P. (2020). *Rollout, policy iteration, and distributed reinforcement learning*. Athena Scientific.
- Bertsekas, D. P., & Tsitsiklis, J. N. (1996). *Neuro-dynamic programming*. Athena Scientific.
- Bertsekas, D. P., Tsitsiklis, J. N., & Wu, C. (1997). Rollout algorithms for combinatorial optimization. *Journal of Heuristics*, 3(3), 245–262. <https://doi.org/10.1023/A:1009635226865>
- Betts, J. T. (2010). *Practical methods for optimal control and estimation using nonlinear programming*. Society for Industrial and Applied Mathematics.
- Bhatnagar, S. (2010, December). An actor-critic algorithm with function approximation for discounted cost constrained Markov decision processes. *Systems & Control Letters*, 59(12), 760–766. <https://doi.org/10.1016/j.sysconle.2010.08.013>
- Borghese, N. A., & Arbib, M. A. (1995). Generation of temporal sequences using local dynamic programming. *Neural Networks*, 8(1), 39–54. [https://doi.org/10.1016/0893-6080\(94\)00053-0](https://doi.org/10.1016/0893-6080(94)00053-0)
- Borkar, V. S. (2005, March). An actor-critic algorithm for constrained Markov decision processes. *Systems & Control Letters*, 54(3), 207–213. <https://doi.org/10.1016/j.sysconle.2004.08.007>
- Boyan, J. A., & Moore, A. W. (1995). Generalization in reinforcement learning: Safely approximating the value function. In *Advances in neural information processing systems* (pp. 369–376). MIT Press.
- Busoniu, L., Babuska, R., De Schutter, B., & Ernst, D. (2010). *Reinforcement learning and dynamic programming using function approximators*. CRC Press.
- Chisci, L., Rossiter, J. A., & Rice, M. J. (1998, April). Stable GPC by dynamic programming. *Systems & Control Letters*, 33(5), 291–300. [https://doi.org/10.1016/S0167-6911\(97\)00116-3](https://doi.org/10.1016/S0167-6911(97)00116-3)
- Clarke, D. W., Mohtadi, C., & Tuffs, P. S. (1987, March). Generalized predictive control—Part I. The basic algorithm. *Automatica*, 23(2), 137–148. [https://doi.org/10.1016/0005-1098\(87\)90087-2](https://doi.org/10.1016/0005-1098(87)90087-2)
- Curry, H. B. (1944). The method of steepest descent for non-linear minimization problems. *Quarterly of Applied Mathematics*, 2(3), 258–261. <https://doi.org/10.1090/qam/1944-02-03>
- Forrest-Barlach, M., & Babcock, S. (1987). Inverse dynamics position control of a compliant manipulator. *IEEE Journal on Robotics and Automation*, 3(1), 75–83. <https://doi.org/10.1109/JRA.1987.1087072>
- Gros, S., & Schild, A. (2017, January). Real-time economic nonlinear model predictive control for wind turbine control. *International Journal of Control*, 90(12), 2799–2812. <https://doi.org/10.1080/00207179.2016.1266514>
- Howard, R. A. (1960). *Dynamic programming and markov processes*. MIT Press.
- Kumar, R. R., & Seywald, H. (1996). Should controls be eliminated while solving optimal control problems via direct methods? *Journal of Guidance, Control, and Dynamics*, 19(2), 418–423. <https://doi.org/10.2514/3.21634>
- Laumond, J., Sekhavat, S., & Lamirault, F. (1998). Guidelines in nonholonomic motion planning for mobile robots. In J. P. Laumond (Ed.), *Robot motion planning and control* (pp. 1–53). Springer.
- Lendaris, G., Cox, C., Saeks, R., & Murray, J. (2002, August 4–7). A radial basis function implementation of the adaptive dynamic programming algorithm. In *The 2002 45th midwest symposium on circuits and systems*, Tulsa, Oklahoma (pp. II–II). IEEE.
- Li, N., Remeikas, C., Xu, Y., Jayasuriya, S., & Ehsani, R. (2015, May). Task assignment and trajectory planning algorithm for a class of cooperative agricultural robots. *Journal of Dynamic Systems, Measurement, and Control*, 137(5), Article 051004. <https://doi.org/10.1115/1.4028849>
- Li, Q., & Xu, Y. (2020, July 1–3). Unmanned aerial vehicle angular velocity control via reinforcement learning in dimension reduced search spaces. In *2020 American control conference*, Denver, CO (pp. 4926–4931). IEEE.
- Liu, D., Huang, Y., Wang, D., & Wei, Q. (2013, May). Neural-network-observer-based optimal control for unknown nonlinear systems using adaptive dynamic programming. *International Journal of Control*, 86(9), 1554–1566. <https://doi.org/10.1080/00207179.2013.790562>
- Luenberger, D. G. (1979). *Introduction to dynamic systems*. John Wiley & Sons.
- Luo, B., Yang, Y., & Liu, D. (2018). Adaptive Q-learning for data-based optimal output regulation with experience replay. *IEEE Transactions on Cybernetics*, 48(12), 3337–3348. <https://doi.org/10.1109/TCYB.2018.2821369>
- Mahony, R., Kumar, V., & Corke, P. (2012). Multirotor aerial vehicles: Modeling, estimation, and control of quadrotor. *IEEE Robotics and Automation Magazine*, 19(3), 20–32. <https://doi.org/10.1109/MRA.2012.2206474>
- Mayne, D. Q., & Michalsha, H. (1988, Dec 7–9). Receding horizon control of nonlinear systems. In *Proceedings of the 27th IEEE conference on decision and control*, Austin, TX (Vol. 1, pp. 464–465). IEEE.
- Michalsha, H., & Mayne, D. Q. (1993, November). Robust receding horizon control of constrained nonlinear systems. *IEEE Transactions on Automatic Control*, 38(11), 1623–1633. <https://doi.org/10.1109/9.262032>
- Miller, W. T., Werbos, P. J., & Sutton, R. S. (1995). *Neural networks for control*. MIT Press.
- Ogata, K. (2010). *Modern control engineering*. Prentice Hall.
- Penrose, R. (1954). A generalized inverse for matrices. In *Mathematical proceedings of the cambridge philosophical society* (Vol. 51, pp. 406–413). Cambridge University Press.
- Pi, C., Hu, K., Cheng, S., & Wu, I. (2020, February). Low-level autonomous control and tracking of quadrotor using reinforcement learning. *Control Engineering Practice*, 95, Article 104222. <https://doi.org/10.1016/j.conengprac.2019.104222>
- Polak, E. (1973). An historical survey of computational methods in optimal control. *SIAM Review*, 15(2), 553–584. <https://doi.org/10.1137/1015071>
- Pontryagin, L. S. (1986). *Mathematical theory of optimal processes*. CRC Press.
- Powell, W. B. (2007). *Approximate dynamic programming: Solving the curses of dimensionality*. John Wiley & Sons.
- Prokhorov, D. V., & Wunsch, D. C. (1997). Adaptive critic designs. *IEEE Transactions on Neural Networks*, 8(5), 997–1007. <https://doi.org/10.1109/TNN.72>
- Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., & Riedmiller, M. (2014, June 22–24). Deterministic policy gradient algorithms. In *Proceedings of the 31st international conference on machine learning*, Beijing, China (pp. 387–395). JMIR.
- Srinivasan, M. V., & Davey, M. (1995). Strategies for active camouflage of motion. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 259(1354), 19–25. <https://doi.org/10.1098/rspb.1995.0004>
- Sundarapandian, V. (2003, January). An invariance principle for discrete-time nonlinear systems. *Applied Mathematics Letters*, 16(1), 85–91. [https://doi.org/10.1016/S0893-9659\(02\)00148-9](https://doi.org/10.1016/S0893-9659(02)00148-9)



- Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1), 9–44. <https://doi.org/10.1007/BF00115009>.
- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT Press.
- Tesauro, G. (1992). Practical issues in temporal difference learning. In *Advances in neural information processing systems* (pp. 259–266). Morgan-Kaufmann.
- Watkins, C. J. C. H. (1989). *Learning from delayed rewards* [Unpublished doctoral dissertation]. UKKing's College.
- Watkins, C. J. C. H., & Dayan, P. (1992, May). Q-learning. *Machine Learning*, 8(3–4), 279–292. <https://doi.org/10.1007/BF00992698>
- Wei, Q., Liu, D., Lin, Q., & Song, R. (2017). Discrete-time optimal control via local policy iteration adaptive dynamic programming. *IEEE Transactions on Cybernetics*, 47(10), 3367–3379. <https://doi.org/10.1109/TCYB.2016.2586082>
- Werbos, P. J. (1982). Applications of advances in nonlinear sensitivity analysis. In F. K. R. F. Drenick (Ed.), *System modeling and optimization* (pp. 762–770). Springer.
- Werbos, P. J. (1989a, June). Back-propagation and neurocontrol: A review and prospectus. In *International 1989 joint conference on neural networks*, Washington, DC(pp. 209–216). IEEE.
- Werbos, P. J. (1989b). Neural networks for control and system identification. In *Proceedings of the 28th IEEE conference on decision and control* (pp. 260–265). IEEE.
- Werbos, P. J. (1990). Consistency of HDP applied to a simple reinforcement learning problem. *Neural Networks*, 3(2), 179–189. [https://doi.org/10.1016/0893-6080\(90\)90088-3](https://doi.org/10.1016/0893-6080(90)90088-3)
- Werbos, P. J. (1992). Approximate dynamic programming for realtime control and neural modelling. In D. A. White & D. A. Sofge (Eds.), *Handbook of intelligent control: Neural, fuzzy and adaptive approaches* (pp. 493–525). Van Nostrand Reinhold.
- Xin, M., & Pan, H. (2009, November). Integrated nonlinear optimal control of spacecraft in proximity operations. *International Journal of Control*, 83(2), 347–363. <https://doi.org/10.1080/00207170903171314>
- Xu, Y., & Basset, G. (2012). Sequential virtual motion camouflage method for nonlinear constrained optimal trajectory control. *Automatica*, 48(7), 1273–1285. <https://doi.org/10.1016/j.automatica.2012.05.017>
- Yaghmaie, F. A., & Braun, D. J. (2019). Reinforcement learning for a class of continuous-time input constrained optimal control problems. *Automatica*, 99, 221–227. <https://doi.org/10.1016/j.automatica.2018.10.038>