WILEY

# Optimal game theoretic solution of the pursuit-evasion intercept problem using on-policy reinforcement learning

**Yusuf Kartal**[1,2] | **Kamesh Subbarao**[2] | **Atilla Dogan**[2] | **Frank Lewis**[1,3]

[1]Automation & Intelligent Systems Division, University of Texas at Arlington Research Institute, Fort Worth, Texas, USA

[2]Mechanical and Aerospace Engineering, University of Texas at Arlington, Arlington, Texas, USA

[3]Electrical Engineering, University of Texas at Arlington, Arlington, Texas, USA

**Correspondence**

Yusuf Kartal, Automation & Intelligent Systems Division, University of Texas at Arlington Research Institute, Forth Worth 76118, TX, USA.
Email: yusuf.kartal@mavs.uta.edu

**Abstract**

This article presents a rigorous formulation for the pursuit-evasion (PE) game when velocity constraints are imposed on agents of the game or players. The game is formulated as an infinite-horizon problem using a non-quadratic functional, then sufficient conditions are derived to prove capture in a finite-time. A novel tracking Hamilton–Jacobi–Isaacs (HJI) equation associated with the non-quadratic value function is employed, which is solved for Nash equilibrium velocity policies for each agent with arbitrary nonlinear dynamics. In contrast to the existing remedies for proof of capture in PE game, the proposed method does not assume players are moving with their maximum velocities and considers the velocity constraints a priori. Attaining the optimal actions requires the solution of HJI equations online and in real-time. We overcome this problem by presenting the on-policy iteration of integral reinforcement learning (IRL) technique. The persistence of excitation for IRL to work is satisfied inherently until capture occurs, at which time the game ends. Furthermore, a nonlinear backstepping control method is proposed to track desired optimal velocity trajectories for players with generalized Newtonian dynamics. Simulation results are provided to show the validity of the proposed methods.

**KEYWORDS**

nonlinear backstepping control, online reinforcement learning, optimal constrained control, pursuit-evasion game, three-dimensional nonlinear systems

## 1 | INTRODUCTION

Inspired by the animal behaviors in hunting scenarios, the pursuit-evasion (PE) games have drawn great attention due to their applicability in areas such as missile guidance,[1] collision avoidance systems[2] and controller designs.[3] The game of this kind is defined as a sub-category of differential game theory and provides the correct framework for the analysis of intercept problem and the choice of optimal policies for the agents involved in the two-player zero-sum (ZS) game.

Isaacs,[4] founder of the differential game theory, initiated the development of strategic policies for both pursuer and evader in a PE problem. In Isaacs,[5] the homicidal chauffeur game was analyzed in detail regarding players' speed and maneuverability capabilities. Bryson[6] introduced optimal feedback laws and demonstrated intercept strategies for players, by using the fixed final-time value function. Lewis et al.[7] made an extension of the Bellman equation, known as the Hamilton-Jacobi-Isaacs (HJI) equations to design $H_\infty$ control, by employing the ZS games solutions. Moreover, works[8-10] deal with linear-quadratic ZS games, in which their objective is to minimize the maximum norm of inputs and states, where the maximum is taken over the unknowns, such as disturbances. Hayoun et al.[11] reveal a set-based computing

method for solving a general class of ZS Stackelberg differential games where the authors come up with a novel class of differential inequalities to get convex outer approximations of backward and forward reachable sets. Bhattacharya et al.[12] worked on a visibility-based PE game when the environment contains a circular obstacle. Furthermore, Li et al.[13] and Liu et al.[14] developed an reinforcement learning (RL) algorithm to learn the Nash equilibrium solution for designing model-free controller by solving the game algebraic Riccati equation forward in time.

Applications of PE games involve the proximate satellite interception guidance strategies. The work[15] studied the intercept problem of satellites where both of the interceptor and target satellite can perform orbital maneuvers with limited thrusts. In the work,[16] the authors analyzed same problem by establishing a local moving coordinate frame and simplifying dynamics of each player to the linear Clohessy-Wiltshire equations. Using the terminal time as a cost function, Gong et al.[17] derived sufficient conditions for capture in the PE problem based on the players' hyper-reachable domain. Note that the common point of these papers is to utilize prescribed terminal time on the construction of the game-theoretical cost function. Jagat et al.[18] proposes quadratic infinite-horizon cost functional for both players but the finite-time capture is not proven mathematically. Instead, simulations are provided to show that capture occurs in a finite-time. Carr et al.[19] employ semi-direct collocation nonlinear programming method to solve optimal actions for agents of the pursuit–evasion game. Authors solve the minimax problem by considering co-state dynamics and boundary conditions simultaneously for the dynamical models.

Standard solution to constrained PE game is to impose external velocity or acceleration constraints. Unfortunately, this leads to discontinuous saturated solutions that are difficult to analyze.[6,7]

Recent works by Hayoun et al.[11] Shaferman et al.[20] and Weintraub et al.[21] focus on the missile-target engagement where the PE problem is formulated as a differential game with an objective of optimizing the linear quadratic cost functional. The work[11] propose bounded maneuverability of the evader to prove the capture in ZS game whereas Weintraub et al.[21] consider an engagement scenario by introducing the defense of a non-maneuverable agent. Further, this work[21] reveals the inclusion of altitude and dynamics in 3-dimensions, which is more realistic for the modeling of aerial engagements.

We sum up the contributions of this article into four categories as:

- First, a backstepping based velocity tracker is developed for PE games where the pursuer and evader both have arbitrary nonlinear dynamics. Taking a priori velocity constraints into account, a novel non-quadratic scalar functional is solved to obtain the smooth optimal velocity policies for each player in contrast to the standard discontinuous solutions.

- Second, with the detailed Lyapunov analysis, sufficient conditions are given for the case where capture must be attained in finite-time.

- The on-policy integral reinforcement learning method is employed to solve the corresponding HJI equation and achieve the game optimal velocity policies for both pursuer and evader.

- Finally, the full rotational dynamics are added to extend the results to full nonlinear dynamical PE systems.

Rest of the article is organized as follows. Section 3 reviews the exponentially stabilizing nonlinear backstepping control method to track given velocity trajectories for a generalized Newtonian system dynamics. Section 4 obtains optimal actions for the players by making use of the Pontryagin's minimum principle and brings analysis of a Nash equilibrium in the PE game. Furthermore, having revealed the sufficient conditions for the asymptotic capture, we prove that PE game ends in a finite-time based on the derived sufficient conditions. Section 5 proposes an on-policy reinforcement learning algorithm for the solution of HJI equation and derives the proof of convergence to the optimal policies. Section 6 closes the backstepping control loop by treating forces and/or moments as finalized inputs to the system and representing the attitude with unit quaternions to overcome the singularity problem of the Euler angles. Finally, the proposed control policies are illustrated via simulation results in Section 7.

## 2 | PROBLEM FORMULATION AND MODEL DESCRIPTION

We study the pursuit-evasion (PE) game for general Newtonian dynamics. A novel approach is given whereby we first design backstepping based velocity controllers for the pursuer and evader that guarantee a Nash solution to the PE game. We use a novel value function that ensures a solution under bounded velocities of the pursuer and the evader. This provides smooth solutions to the bounded velocity PE game in contrast to standard discontinuous solutions.[6] We conduct

a Nash equilibrium analysis for a game of this kind. Further, we seek to obtain sufficient conditions for global exponential stability of the origin (equilibrium) using a rigorous Lyapunov analysis. Finally, we seek to derive conditions for a final time capture, and provide an upper bound on the time of capture.

The generalized translational and rotational dynamics for the pursuer and the evader can be modeled in their respective body frames of reference as

$$m^i \dot{v}_B^i = m^i S(w_B^i) v_B^i + N^i(\eta^i) f_g^i + f_B^i, \tag{1}$$

$$I_B^i \dot{w}_B^i = S(w_B^i) I_B^i w_B^i + \tau_B^i \tag{2}$$

where the superscript $i \in \{p, e\}$ with $p$ denoting the pursuer and $e$ denoting the evader respectively. Here, $v_B^i \in \mathbb{R}^3$, $w_B^i \in \mathbb{R}^3$ are the translational and angular velocities respectively, and $f_B^i \in \mathbb{R}^3$, $\tau_B^i \in \mathbb{R}^3$ are the control forces and moments respectively, in the body fixed reference frame. Further, $I_B^i \in \mathbb{R}^{3\times3}$ is the constant nonsingular inertia matrix defined in the body frame and $m^i$ is a scalar quantity that denotes the mass of players' rigid bodies. In addition, $f_g^i = \begin{bmatrix} 0 & 0 & m^i g \end{bmatrix}^T$ is the gravitational force vector whose components are written in the Inertial frame. $S(w_B^i) \in \mathbb{R}^{3\times3}$ represents a skew-symmetric matrix form of the vector $w_B^i$. Moreover, $m^i S(w_B^i) v_B^i$ and $S(w_B^i) I_B^i w_B^i$ are due to the derivative of the body referenced linear and angular momentum of the vehicles relative to the Inertial frame. $N^i(\eta^i) \in \mathbb{R}^{3\times3}$ is a rotation matrix from Inertial to body frame with the argument of Euler angle vector $\eta^i \in \mathbb{R}^3$ (see Equation (46) in Section 6). Later in Section 6, we will call this Inertial frame as earth frame and give detailed explanation for the rotation matrix.

## 3 | DEVELOPING VELOCITY TRACKER USING BACKSTEPPING CONTROL METHOD

In this section, we present an exponentially stabilizing backstepping control method to track given velocity trajectories. This velocity tracker is developed in this section, which uses only the translational dynamics (1). In Section 4, the velocity tracker is extended for PE games based on the translational dynamics (1) for both pursuer and evader. Then, in Section 6 we also consider rotational dynamics (2) to obtain general controllers for both velocity and attitude for pursuer and evader.

Note, we first derive the required velocity tracking control laws in the Inertial frame, and then subsequently Section 6 shows how they are realized using the dynamics in (1) and (2).

Using standard techniques,[22] the translational dynamics (1) is represented in the Inertial frame as,

$$m^i \dot{v}^i = f^i + f_g^i \tag{3}$$

where $v^i \in \mathbb{R}^3$ is the velocity vector and $f^i = N^{i^T}(\eta^i) f_B^i$ is the control force, in the Inertial frame $i \in \{p, e\}$. Here $N^i(\eta^i)$ is a rotation matrix given in (46) (Section 6) that depends on the Euler angles, $\eta^i = \begin{bmatrix} \psi^i & \theta^i & \varphi^i \end{bmatrix}^T$, generated by the rotational system (2). Therefore, $f^i$ cannot be directly controlled, but depends on the rotational dynamics (2). See Section 6 for elaboration.

Therefore, backstepping must be used to determine the desired $f^i$ that must be generated by (2). Introducing a desired virtual force $f_d^i$, to the system dynamics (3) we obtain

$$m^i \dot{v}^i = f_d^i + f_g^i + \tilde{f}^i \tag{4}$$

where $\tilde{f}^i = f^i - f_d^i$ is the difference of control and desired forces of the Newtonian system in 3-D.

Define velocity error as

$$\delta_v^i = v_d^i - v^i \tag{5}$$

where $v_d^i \in \mathbb{R}^3$ is the desired velocity designed for pursuer $v_d^p$ and evader $v_d^e$ in the next section. Take the derivative of (5) and substitute in (4) to obtain closed-loop velocity error dynamics as

$$m^i \dot{\delta}_v^i = m^i \dot{v}_d^i - f_g^i - \tilde{f}^i - f_d^i. \tag{6}$$

Then select ideal desired force as

$$\boldsymbol{f}_d^i = m^i \dot{\boldsymbol{v}}_d^i - \boldsymbol{f}_g^i + m^i \boldsymbol{K}^i \boldsymbol{\delta}_v^i \tag{7}$$

where $\boldsymbol{K}^i \in \mathbb{R}^{n \times n}$ is a positive-definite matrix. Substituting (7) in (6) yields

$$\dot{\boldsymbol{\delta}}_v^i = -\boldsymbol{K}^i \boldsymbol{\delta}_v^i - \frac{\widetilde{\boldsymbol{f}}^i}{m^i}. \tag{8}$$

This enables us to derive exponential stability of the origin, as long as an admissible $\boldsymbol{f}_d^i$ exists. In Section 6 we consider the rotational dynamics and show how to design the control force, $\dot{\boldsymbol{f}}^i$ and hence $\dot{\boldsymbol{f}}_B^i$ in (1) and (3) respectively, to make $\widetilde{\boldsymbol{f}}^i \to \boldsymbol{0}$.[23] Then (8) shows that $\boldsymbol{\delta}_v^i \to \boldsymbol{0}$ exponentially.

*Remark* 1. Tracking the vector quantity $\boldsymbol{f}_d^i$ in (7) not only guarantees exponential stability of the equilibrium of (6) but also gives the desired attitude of the Newtonian system (3) so that it is aligned with the direction of $\boldsymbol{f}_d^i$.

The next section deals with the derivation of optimal velocity trajectories $\boldsymbol{v}_d^i$ for pursuer and evader, employed in (5). The design of the desired ideal forces $\boldsymbol{f}_d^p, \boldsymbol{f}_d^e$ is treated in Section 6.

## 4 | OPTIMAL GAME THEORETIC VELOCITY GENERATION FOR PURSUIT-EVASION GAME

In this main section, we first propose a formulation of PE game and derive the optimal bounded desired velocity trajectories $\boldsymbol{v}_d^p, \boldsymbol{v}_d^e$ in (5) for the players. Second, we conduct a Nash equilibrium analysis for the game and derive sufficient conditions for global exponential stability of the origin by rigorous Lyapunov analysis. Finally, conditions for finite-time capture and its upper bound are given.

### 4.1 | Pursuit-evasion game formulation

Assuming the players are governed by the velocity error dynamics (8), this section presents various definitions to develop the game-theoretically optimal solution of the PE game satisfying *velocity constraints* on the players. To simplify the notation, define desired velocity in (5) for the pursuer $\boldsymbol{v}^p = \boldsymbol{v}_d^p$ and the evader $\boldsymbol{v}^e = \boldsymbol{v}_d^e$.

The following kinematic expressions enable us to derive desired velocities and thereby the forces (7) for pursuer and evader

$$\dot{\boldsymbol{\xi}}^p = \boldsymbol{v}^p$$
$$\dot{\boldsymbol{\xi}}^e = \boldsymbol{v}^e \tag{9}$$

where $\boldsymbol{\xi}^p \in \mathbb{R}^3$ and $\boldsymbol{\xi}^e \in \mathbb{R}^3$ denote the 3-dimensional position vectors $(x, y, z)$ of pursuer and evader respectively, which are defined with respect to Inertial frame. Hence $\boldsymbol{v}^p \in \mathbb{R}^3$ and $\boldsymbol{v}^e \in \mathbb{R}^3$ are desired velocity vectors of the pursuer and evader respectively. Note that (3) employs the translational velocity in the PE game. This allows analysis of ZS game for general nonlinear systems in Section 6.

Now, consider the following formulation for the zero-sum (ZS) PE game. Let the evader have an objective of maximizing the relative distance $\boldsymbol{\delta} \in \mathbb{R}^3$, defined as

$$\boldsymbol{\delta} = \boldsymbol{\xi}^p - \boldsymbol{\xi}^e, \tag{10}$$

whereas the pursuer tries to minimize (10). Moreover, let the velocities of both pursuer and evader be bounded by scalars $|v_j^p| \leq \lambda^p; |v_j^e| \leq \lambda^e \forall j = 1, \dots, n$. To satisfy these constraints, the value functional is defined as

$$V^{\pi^p, \pi^e}(\boldsymbol{\delta}) = \int_t^\infty \{\boldsymbol{\delta}^T \boldsymbol{Q} \boldsymbol{\delta} + U(\pi^p(\boldsymbol{\delta})) - U(\pi^e(\boldsymbol{\delta}))\} d\tau \tag{11}$$

where $Q \in \mathbb{R}^{nxn}$ is a positive-definite matrix, $\pi^p(.)$ and $\pi^e(.)$ stand for the policies of pursuer and evader respectively in ZS game such that

$$\pi^p(\delta) \triangleq v^p$$
$$\pi^e(\delta) \triangleq v^e. \tag{12}$$

Moreover, $U(v^i)$ (for $i$ is either $p$ or $e$) is a generalized non-quadratic scalar functional,[24] which ensures bounded velocities given by

$$U(v^i) = 2 \int_0^{v^i} (\alpha^{-1}(u^i/\lambda^i))^T R^i du^i,$$

$$\alpha^{-1}(u^i/\lambda^i) = \left[ \alpha^{-1}(u_1^i/\lambda^i) \; \ldots \; \alpha^{-1}(u_n^i/\lambda^i) \right]^T,$$

$$u^i = \left[ u_1^i \; \ldots \; u_n^i \right]^T, v^i = \left[ v_1^i \; \ldots \; v_n^i \right]^T \tag{13}$$

where $R^i \in \mathbb{R}^{nxn}$ is a symmetric positive-definite matrix and $\alpha(.)$ is a bounded one-to-one smooth function that is, it belongs to $C^\ell, \ell \geq 1$. This is a monotonic odd function with its first derivative bounded by a constant. An example of $\alpha(.)$ is tanh(.) and throughout this article, we use tanh(.), which constrains the velocity to remain within predefined limits that is, $|v_j^i| \leq \lambda, \forall j = 1, \ldots, n$ and $\forall i = p, e$. In ZS PE games, $R^i$ plays a key role by restricting the rate of change of optimal velocities and hence constrains the accelerations of the each player.

The differential equivalent of (11) is the ZS game Bellman equation. Using (9), (10) and Leibniz's formula, the Bellman equation is obtained as

$$H(\delta, \nabla V, v^p, v^e) \equiv \delta^T Q \delta + U(v^p) - U(v^e) + \nabla V^T \dot{\delta}$$
$$\equiv \delta^T Q \delta + U(v^p) - U(v^e) + \nabla V^T (v^p - v^e) = 0 \tag{14}$$

where $\nabla V = \partial V^{\pi^p, \pi^e}/\partial \delta \in \mathbb{R}^n$ is the gradient of value function (11), and $H(.)$ is the Hamiltonian.

To find the optimal policies $\pi^{i*}(\delta) = v^{i*}$ (for $i = p, e$) of players in the game, check stationarity conditions $\partial H/\partial v^p = 0$ and $\partial H/\partial v^e = 0$. For the pursuer, applying Pontryagin's minimum principle to (14) yields

$$\frac{\partial H}{\partial v^p} \equiv \frac{\partial U(v^p)}{\partial v^p} + \frac{\partial}{\partial v^p} \{ \nabla V^T (v^p - v^e) \}. \tag{15}$$

Evaluating the derivatives at the right-hand-side of (15) using Leibniz's formula, and checking the stationarity condition $\partial H/\partial v^p = 0$ yields

$$2 \left( \tanh^{-1} \left( \frac{v^{p*}}{\lambda^p} \right) \right)^T R^p = -\nabla V^{*T}. \tag{16}$$

Then, the optimal policy for the pursuer using the definition (12) is obtained as

$$\pi^{p*}(\delta) \triangleq v^{p*} = -\lambda^p \tanh \left( \frac{1}{2} (R^p)^{-1} \nabla V^* \right). \tag{17}$$

This velocity control bounded as required.

Likewise, one can follow the same steps to derive bounded optimal velocity policy for the evader as

$$\pi^{e*}(\delta) \triangleq v^{e*} = -\lambda^e \tanh \left( \frac{1}{2} (R^e)^{-1} \nabla V^* \right). \tag{18}$$

Let $V^*$ be the optimal value of (11) with the policies given in (17) and (18). Then Hamilton-Jacobi-Isaacs (HJI) equation is obtained as

$$H(\delta, \nabla V^*, v^{p*}, v^{e*}) \equiv \delta^T Q \delta + U(v^{p*}) - U(v^{e*}) + \nabla V^{*T} (v^{p*} - v^{e*}) = 0. \tag{19}$$

Note that the positive and negative definiteness of Hessians, $\partial^2 H/\partial v^{p2} > 0$ and $\partial^2 H/\partial v^{e2} < 0$, indeed show that pursuer's optimal policy aims to minimize the Hamiltonian (14) whereas evader's aims to maximize. Therefore, $(v^{p*}, v^{e*})$ is the game-theoretic saddle point. Furthermore, this is a Nash equilibrium since the game is of type ZS and (11) is separable.[6] Rigorous analysis of this is shown in Theorem 1.

## 4.2 | Proof of Nash equilibrium

In this section, we derive the value of PE game at Nash equilibrium. The following lemmas and corollary are necessary steps to prove that the Nash equilibrium is reached with policies (17) and (18).

**Lemma 1.** *Let $V^{\pi^p,\pi^e}(\delta)$ be the corresponding solution of the Hamiltonian (14). Then following equality holds*

$$
\begin{aligned}
H(\delta, \nabla V, v^p, v^e) = H(\delta, \nabla V, v^{p*}, v^{e*}) + \nabla V^T((v^p - v^{p*}) \\
+ (v^{e*} - v^e)) + U(v^p) - U(v^{p*}) + U(v^{e*}) - U(v^e).
\end{aligned}
\tag{20}
$$

*Proof.* Adding and subtracting the terms $U(v^{p*})$, $U(v^{e*})$, $\nabla V^T v^{p*}$, and $\nabla V^T v^{e*}$ to Hamiltonian (14) yields

$$
\begin{aligned}
H(\delta, \nabla V, v^p, v^e) = \delta^T Q \delta + \nabla V^T(v^{p*} - v^{e*}) U(v^{p*}) \\
- U(v^{e*}) + \nabla V^T((v^p - v^{p*}) + (v^{e*} - v^e)) \\
+ U(v^p) - U(v^{p*}) + U(v^{e*}) - U(v^e),
\end{aligned}
\tag{21}
$$

which completes the proof. ∎

**Lemma 2.** *Let $V^{\pi^p,\pi^e}(\delta)$ be the corresponding solution of the Hamiltonian (14) and define $V(\delta(t_0))$ as the initial value of the game. Then following equality holds*

$$
V^{\pi^p,\pi^e}(\delta(t_0)) = \int_{t_0}^{\infty} H(\delta, \nabla V, v^p, v^e) d\tau + V(\delta(t_0)).
\tag{22}
$$

*Proof.* Assume that capture occurs in the interval $t \in [t_0, \infty]$, which implies $\lim_{t \to \infty} V^{\pi^p,\pi^e}(\delta(t)) = 0$. Then adding zero to (11) yields

$$
\begin{aligned}
V^{\pi^p,\pi^e}(\delta(t_0)) &= \int_{t_0}^{\infty} \{\delta^T Q \delta + U(\pi^p(\delta)) - U(\pi^e(\delta))\} d\tau + \int_{t_0}^{\infty} \dot{V}^{\pi^p,\pi^e} d\tau + V(\delta(t_0)) \\
&= \int_{t_0}^{\infty} \{\delta^T Q \delta + U(v^p) - U(v^e)\} d\tau + \int_{t_0}^{\infty} \nabla V^T(v^p - v^e) d\tau + V(\delta(t_0)) \\
&= \int_{t_0}^{\infty} H(\delta, \nabla V, v^p, v^e) d\tau + V(\delta(t_0)).
\end{aligned}
\tag{23}
$$

This completes the proof. ∎

The next corollary extends the fact given in Lemma 1.

**Corollary 1.** *Suppose $V^*$ satisfies the HJI Equation 19. Then $H(\delta, \nabla V^*, v^{p*}, v^{e*}) = 0$ and (20) becomes*

$$
H(\delta, \nabla V^*, v^p, v^e) = \nabla V^{*T}((v^p - v^{p*}) + (v^{e*} - v^e)) + U(v^p) - U(v^{p*}) + U(v^{e*}) - U(v^e).
\tag{24}
$$

The next theorem derives the optimal value of the ZS game and proves Nash equilibrium reached.

**Theorem 1.** *Consider kinematic expressions for the players (9) with the value function given in (11). Let $V^*$ be a positive definite smooth solution of HJI Equation 19. Then, $(v^{p*}, v^{e*})$ given by (17), (18) is the Nash equilibrium and $V^*(\delta(t_0))$ is the value of PE game.*

*Proof.* Using the facts given in Lemma 2 and Corollary 1, (23) becomes

$$V^{\pi^p,\pi^e}(\boldsymbol{\delta}(t)) = \int_t^\infty \{\nabla V^{*T}((\boldsymbol{v^p} - \boldsymbol{v^{p*}}) + (\boldsymbol{v^{e*}} - \boldsymbol{v^e}))$$
$$+ U(\boldsymbol{v^p}) - U(\boldsymbol{v^{p*}}) + U(\boldsymbol{v^{e*}}) - U(\boldsymbol{v^e})\}d\tau + V^*(\boldsymbol{\delta}(t_0)). \tag{25}$$

To prove $(\boldsymbol{v^{p*}}, \boldsymbol{v^{e*}})$ is the Nash equilibrium of the game, we need to show that when the pursuer adopts policy given in (17), the best action for the evader to maximize the value function (11) is $\boldsymbol{v^{e*}}$. Likewise, when the evader adopts policy given in (18), the best action for the pursuer to minimize the value function (11) is $\boldsymbol{v^{p*}}$ that is,

$$V^{\pi^{p*},\pi^e}(\boldsymbol{\delta}(t)) \leq V^{\pi^{p*},\pi^{e*}}(\boldsymbol{\delta}(t)) \leq V^{\pi^p,\pi^{e*}}(\boldsymbol{\delta}(t)). \tag{26}$$

Note that $V^{\pi^{p*},\pi^{e*}}(\boldsymbol{\delta}(t)) = V^*(\boldsymbol{\delta}(0))$ and call the integral term in (25) as $\beta(V^{\pi^p,\pi^e})$. Now we need to show $\beta(V^{\pi^{p*},\pi^e}) \leq 0$ and $\beta(V^{\pi^p,\pi^{e*}}) \geq 0$ so that (26) holds. Then using (13), (17), (18), and (25) we obtain

$$\beta(V^{\pi^{p*},\pi^e}) = \int_t^\infty \{\nabla V^{*T}(\boldsymbol{v^{e*}} - \boldsymbol{v^e}) + U(\boldsymbol{v^{e*}}) - U(\boldsymbol{v^e})\}d\tau$$
$$= \int_t^\infty \left\{ -2(\tanh^{-1}(\boldsymbol{v^{e*}}/\lambda^e))^T \boldsymbol{R^e}(\boldsymbol{v^{e*}} - \boldsymbol{v^e}) + 2\int_{\boldsymbol{v^e}}^{\boldsymbol{v^{e*}}} (\tanh^{-1}(\boldsymbol{u}/\lambda^e))^T \boldsymbol{R^e}d\boldsymbol{u} \right\} d\tau. \tag{27}$$

Now define $\phi^T(.) = \tanh^{-1}(.)$ and note that $\phi^T(.)$ is monotonically increasing function in the interval $[-\lambda^e, \lambda^e]$. To complete the proof, first assume that $\boldsymbol{v^{e*}} \geq \boldsymbol{v^e}$ and apply integral mean value theorem on (27)

$$\beta(V^{\pi^{p*},\pi^e}) = \int_t^\infty \left\{ -2\phi(\boldsymbol{v^{e*}}/\lambda^e)\boldsymbol{R^e}(\boldsymbol{v^{e*}} - \boldsymbol{v^e}) + 2\int_{\boldsymbol{v^e}}^{\boldsymbol{v^{e*}}} \phi(\boldsymbol{u}/\lambda^e)\boldsymbol{R^e}d\boldsymbol{u} \right\} d\tau$$
$$\leq \int_t^\infty \{-2\phi(\boldsymbol{v^{e*}}/\lambda^e)\boldsymbol{R^e}(\boldsymbol{v^{e*}} - \boldsymbol{v^e}) + 2\phi(\boldsymbol{v^{e*}}/\lambda^e)\boldsymbol{R^e}(\boldsymbol{v^{e*}} - \boldsymbol{v^e})\}d\tau = 0. \tag{28}$$

Then assume that $\boldsymbol{v^{e*}} < \boldsymbol{v^e}$ and again apply integral mean value theorem on (27)

$$\beta(V^{\pi^{p*},\pi^e}) = \int_t^\infty \left\{ 2\phi(\boldsymbol{v^{e*}}/\lambda^e)\boldsymbol{R^e}(\boldsymbol{v^e} - \boldsymbol{v^{e*}}) - 2\int_{\boldsymbol{v^{e*}}}^{\boldsymbol{v^e}} \phi(\boldsymbol{u}/\lambda^e)\boldsymbol{R^e}d\boldsymbol{u} \right\} d\tau$$
$$\leq \int_t^\infty \{2\phi(\boldsymbol{v^{e*}}/\lambda^e)\boldsymbol{R^e}(\boldsymbol{v^e} - \boldsymbol{v^{e*}}) - 2\phi(\boldsymbol{v^{e*}}/\lambda^e)\boldsymbol{R^e}(\boldsymbol{v^e} - \boldsymbol{v^{e*}})\}d\tau = 0, \tag{29}$$

which shows that $\beta(V^{\pi^{p*},\pi^e}) \leq 0$. The same procedure can be performed to show $\beta(V^{\pi^p,\pi^{e*}}) \geq 0$. Then the inequality given in (26) is verified, which implies that $(\boldsymbol{v^{p*}}, \boldsymbol{v^{e*}})$ is the Nash equilibrium and $V^*(\boldsymbol{\delta}(t_0))$ is the value of the PE game. ∎

## 4.3 | Stability and finite-time capture analysis

This section first reveals the sufficient conditions for the asymptotic capture of the evader by the pursuer. Then, by making use of these conditions, derives the globally exponential stability of the origin. Finally it is shown that under certain conditions, finite-time capture is ensured.

Before developing analysis for the asymptotic capture, let $\boldsymbol{R^i}$ in (17) and (18) be a diagonal matrix with elements of $r_j^i > 0, \forall j \in \{1, 2, 3\}$ and $\forall i = p, e$. This enables us to simplify the analysis that will be developed in the rest of the article. Employing this assumption, the next theorem shows the sufficient conditions for the asymptotic capture in ZS PE games.

**Theorem 2.** *Consider kinematic expressions for the players (9) with the value function given in (11). Then the equilibrium of tracking error dynamics $\dot{\boldsymbol{\delta}} = \boldsymbol{v^{p*}} - \boldsymbol{v^{e*}}$, is asymptotically stable point with candidate Lyapunov function $L(\boldsymbol{\delta}) = V^{\pi^{p*},\pi^{e*}}(\boldsymbol{\delta})$. The sufficient conditions for asymptotic capture are $\lambda^p > \lambda^e$ and $r_{e_i} \geq r_{p_i} \forall i = 1, \dots, n$.*

*Proof.* Since $V^{\pi^{p*},\pi^{e*}}(\boldsymbol{\delta})$ does not depend on the time explicitly, equality $\dot{L}(\boldsymbol{\delta}) = \nabla L^T \dot{\boldsymbol{\delta}}$ holds. By (14), derivative of the Lyapunov function, $\dot{L}(\boldsymbol{\delta})$ is obtained as

$$\dot{L}(\boldsymbol{\delta}) = -\boldsymbol{\delta}^T \boldsymbol{Q}\boldsymbol{\delta} - U(\boldsymbol{v}^{p*}) + U(\boldsymbol{v}^{e*}). \tag{30}$$

Assumption of the equality $\boldsymbol{R}^p = \boldsymbol{R}^e$, implies the pursuer and evader are moving in the same direction by (17) and (18). For intercept, the position of the pursuer and evader must be equal. To meet this criteria, we propose $\lambda^p > \lambda^e$ so that asymptotic capture occurs as the row elements of optimal actions satisfy $|v^{p*}|_i > |v^{e*}|_i \; \forall i = 1, \ldots, n$. Furthermore, taking (30) into account, the condition of $\boldsymbol{R}^p = \boldsymbol{R}^e$ is relaxed as $\boldsymbol{R}^e \geq \boldsymbol{R}^p$ since proposition $\lambda^p > \lambda^e$ implies $U(\boldsymbol{v}^{p*}) \geq U(\boldsymbol{v}^{e*})$, $\dot{L}(\boldsymbol{\delta})$ becomes strictly negative definite. Then sufficient conditions for the asymptotic capture is proved to be $\lambda^p > \lambda^e$ and $r_{e_i} \geq r_{p_i} \forall i = 1, \ldots, n$. ∎

*Remark* 2. Aysmptotic capture in Theorem 2 can be strengthened to finite-time capture with the assumption that players involved in the game satisfy the sufficient conditions derived in the proof of Theorem 2. See Lemma 3.

Following theorem extends the Theorem 2 to exponential stability of the origin.

**Theorem 3.** *Consider sufficient conditions and Lyapunov function, $L(\boldsymbol{\delta})$ given in Theorem 2. Then, there exists positive scalars $c_1$, $c_2$, and $\epsilon$, which satisfies*

$$c_1 \|\boldsymbol{\delta}\|_2^2 \leq L(\boldsymbol{\delta}) \leq c_2 \|\boldsymbol{\delta}\|_2^2$$
$$\dot{L}(\boldsymbol{\delta}) \leq -\epsilon L(\boldsymbol{\delta}), \tag{31}$$

*which implies that the origin is an exponentially stable equilibrium. Furthermore, radially unboundedness of the $L(\boldsymbol{\delta})$ implies the globally exponentially stability of the origin,[25] which is an essential result as the initial positional offset between the pursuer and evader should not be problem to prove the capture in PE game.*

*Proof.* The inequality $U(\boldsymbol{v}^{p*}) \geq U(\boldsymbol{v}^{e*})$ by Thoerem 2 and the strict convexity of $U(\boldsymbol{v}_i)$ (for $i = p, e$), imply the existence of positive scalars $c_1$ and $c_2$.[26] Now, define convex function $U_s(\boldsymbol{\delta})$ that satisfies the inequality $U_s(\boldsymbol{\delta}) \leq U(\boldsymbol{v}^{p*}) - U(\boldsymbol{v}^{e*})$. Using this and (30), the following inequality is derived as

$$\dot{L}(\boldsymbol{\delta}) \leq -\boldsymbol{\delta}^T \boldsymbol{Q}\boldsymbol{\delta} - U_s(\boldsymbol{\delta}). \tag{32}$$

Substituting (32) in (11) with the optimal policies (17) and (18), results in

$$L(\boldsymbol{\delta}) \leq \int_t^\infty \{\boldsymbol{\delta}^T \boldsymbol{Q}\boldsymbol{\delta} + U_s(\boldsymbol{\delta})\} d\tau, \tag{33}$$

which stands for the proof of $\dot{L}(\boldsymbol{\delta}) \leq -\epsilon L(\boldsymbol{\delta})$ for sufficiently small $\epsilon$, which completes the proof. ∎

Notice that PE game given in Section 4 is formulated by treating the players as unit masses since the kinematic expressions (9) is employed in the value function (11). In Section 6, we consider full nonlinear dynamics (1), (2). Now, consider the volume of pursuer and evader in 3-dimensional space and let the pursuer and evader have a sphere of collision with radius $r^p$ and $r^e$ respectively, as illustrated in Figure 1. Then capture occurs when the distance between the center of masses of players is less than $r^p + r^e$. With this in mind, the next main lemma proves that the capture of evader by pursuer indeed occurs in finite-time in PE game.

**Lemma 3.** *There exists an upper-bound for the capture time in PE game when the conditions $\lambda^p > \lambda^e$ and $r_{e_i} \geq r_{p_i} \forall i = 1, \ldots, n$ derived in Theorem 2 are satisfied. This also implies that the PE game ends in finite-time as required.*

*Proof.* The globally exponentially stability of the origin derived in Theorem 3 implies the equation of positional offset between the players is in the form of

$$\|\boldsymbol{\delta}(t)\|_2 \leq b_1 \|\boldsymbol{\delta}(t_0)\|_2 e^{-b_2(t-t_0)} \quad \forall t > t_0 \tag{34}$$
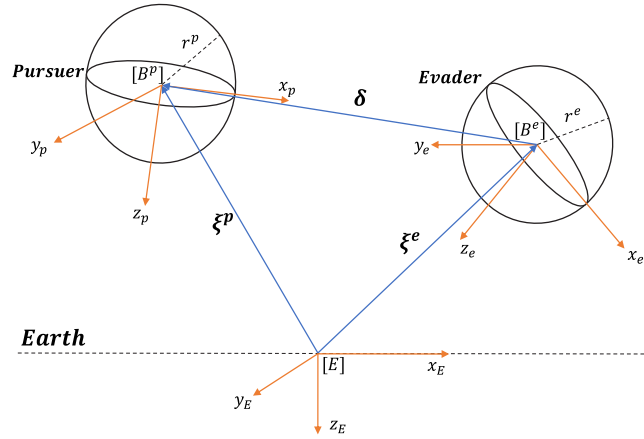
**FIGURE 1** Sphere of collisions for players and their frames in 3-dimensions used for finite-time capture analysis [Colour figure can be viewed at wileyonlinelibrary.com]

where $b_i$ is a positive scalar $\forall i = 1, 2$. Then the upper bound for capture time $t_c$ is derived as

$$t_c \leq t_0 + \frac{1}{b_2} \log \left( \frac{b_1 \|\delta(t_0)\|_2}{r^e + r^p} \right) \tag{35}$$

where $\log(.)$ is a natural logarithm function and this completes the proof. ∎

*Remark* 3. It is seen that for finite-time capture, the velocity bound $\lambda^p$ for the pursuer must be greater then the velocity bound $\lambda^e$ on the evader. Moreover, the sufficient condition on weights (13) is found as $r_{e_i} \geq r_{p_i} \forall i = 1, \ldots, n$. Note that capture time is also studied for multi-agent systems in the work[27] by assuming the players are using their maximum efforts. In Lemma 3, we showed that capture time is upper bounded under certain conditions even the players are not using their maximum efforts.

# 5 | ONLINE SOLUTION OF HJI EQUATION USING INTEGRAL REINFORCEMENT LEARNING (IRL)

The PE game formulation in Section 4 requires the generation of velocity set-points online and in real-time for both agents of the game. With this in mind, we employ the following synchronous IRL algorithm[28] to solve the HJI Equation 19 in real-time and hence, reach the Nash equilibrium velocity policies $(\nu^{p*}, \nu^{e*})$ online by observing measured data. In the work [28], it is emphasized that persistence of excitation condition must be satisfied so that the IRL algorithm convergences. This is achieved in most applications[28,29] by adding small probing noise. In our case, the persistence of excitation for IRL to work is satisfied inherently until capture occurs, at which time the game ends.

The tracking HJI Equation 19 is nonlinear in the value function gradient $\nabla V^*$, and non-quadratic partial differential equation that is extremely difficult to solve. However, (36) can be solved for the value function and its gradient by collecting position data over some interval $[t, t + T]$. Therefore, finding the value of game optimal velocity policies by solving (36) is easier than solving (19). This is the motivation of introducing an iterative algorithm for approximating the tracking HJI solution, which is necessary to evaluate game optimal velocity policy for pursuer (17), and evader (18).

## 5.1 | Policy iteration solution for PE game

In this section, we present a policy iteration algorithm that avoids solution of (19) and also does not require knowledge of the system dynamics. The following lemma enables us to recognize IRL form of the value function (11).

**Lemma 4.** *Let $V^{\pi^p,\pi^e}(\delta)$ be the corresponding solution of the Bellman Equation 14. Then, the value function (11), can be written in the IRL form as*

$$V^{\pi^p,\pi^e}(\delta(t)) = \int_t^{t+T} \{\delta^T Q\delta + U(\pi^p(\delta)) - U(\pi^e(\delta))\}d\tau + V^{\pi^p,\pi^e}(\delta(t+T)). \tag{36}$$

*Proof.* The equality $\dot{V}^{\pi^p,\pi^e} = -\delta^T Q\delta - U(\pi^p(\delta)) + U(\pi^e(\delta))$ holds by the differentiation of ZS game Bellman Equation 14. Then integrating both sides from $t$ to $t+T$, results in

$$\int_t^{t+T} \dot{V}d\tau = -\int_t^{t+T} \{\delta^T Q\delta + U(\pi^p(\delta)) - U(\pi^e(\delta))\}d\tau, \tag{37}$$

which verifies (36). ∎

The online policy-iteration Algorithm 1 performs a sequence of four-step iterations to find the optimal control policies for players. Notice that these policies stand for the optimal desired velocities, which are employed in (5). Furthermore, they are also Nash equilibrium velocity policies by Theorem 1.

---

**Algorithm 1.** Online policy-iteration algorithm

---

1. Select any policy $\pi_0^p$ and $\pi_0^e$ for the players
2. Policy evaluation

$$V^{\pi_j^p,\pi_j^e}(\delta(t)) = \int_t^{t+T} \{\delta^T Q\delta + U(\pi_j^p(\delta)) - U(\pi_j^e(\delta)\}d\tau + V^{\pi_j^p,\pi_j^e}(\delta(t+T)). \tag{38}$$

3. Policy improvement

$$\pi_{j+1}^p(\delta) = -\lambda^p \tanh\left(\frac{1}{2}(R^p)^{-1}\nabla V^{\pi_j^p,\pi_j^e}\right),$$
$$\pi_{j+1}^e(\delta) = -\lambda^e \tanh\left(\frac{1}{2}(R^e)^{-1}\nabla V^{\pi_j^p,\pi_j^e}\right). \tag{39}$$

4. On convergence stop; else go to step 2. ☐

---

Notice that the position data of each player is collected through each iteration over the period $T$. The proof of convergence of Algorithm 1 to the optimal policies is shown in the following theorem.

**Theorem 4.** *Using the temporal difference (TD) learning method, Algorithm 1 converges to the Nash value $V^*(\delta(t_0))$ and Nash equilibrium policies $(\pi^{p*}, \pi^{e*})$, which optimizes velocity trajectories for the players in a game theoretic manner.*

*Proof.* First, evaluate the value function $V^{\pi^{pj},\pi^{ej}}(\delta(t))$, which solves the (38) by TD method. Then by Theorem 1, Isaacs' condition is derived as

$$H(\delta, \nabla V, v^{p*}, v^e) \le H(\delta, \nabla V, v^{p*}, v^{e*}) \le H(\delta, \nabla V, v^p, v^{e*}). \tag{40}$$

Noting $\beta(V^{\pi^{p*},\pi^e}) \le 0$ and $\beta(V^{\pi^p,\pi^{e*}}) \ge 0$ is proved in the Theorem 1, the uniform convergence of Algorithm 1 immediately follows from Dini's theorem as reinforcement $H(\delta, \nabla V^j, v^p, v^e)$ converges to $H(\delta, \nabla V^*, v^{p*}, v^{e*}) = 0$ by Corollary 1. Moreover, due to the uniqueness of the value function (11), it follows that $\lim_{j\to\infty} V^{\pi_j^p,\pi_j^e}(\delta(t)) = V^*(\delta(t_0))$. ∎

## 5.2 | Value function approximation to find game optimal pursuer and evader velocity policies

This section presents a critic neural network structure for policy-evaluation step in Algorithm 1.

*Remark* 4. The IRL method given in Algorithm 1 requires the value function approximation (VFA), which can be achieved in a least-squares sense that is also known as single hidden layer critic Neural Network (NN). We employ this technique as in Reference 28 that guarantees the successive least-squares iterations converge to the optimal value function of the HJI Equation 19, and hence $\nabla V^*$.

*Remark* 5. Note that the pair $(\boldsymbol{\nu}^{p*}, \boldsymbol{\nu}^{e*})$ stands for the Nash equilibrium by Theorem 1, thereby the Algorithm 1 converges to optimal actions for both players. Unlike the works[29,30] that use the IRL technique to reach min or max point of the value functional, we employ this technique to converge game theoretic saddle point by using the Isaacs' condition derived in Theorem 4. In addition, the system dynamics (1) does not appear in the value functional, which implies that we do not need to implement actor NN[30,31] and the solution of HJI (19) can be obtained by using only critic NN, see Reference 32.

By Remarks 4 and 5, we approximate the game optimal value functional in step 2 of Algorithm 1 using Weierstrass approximator such that

$$\hat{V}(\delta) = \hat{\boldsymbol{W}}^T \boldsymbol{\Phi}(\delta),$$
$$\nabla \hat{V} = \nabla \boldsymbol{\Phi}(\delta)^T \hat{\boldsymbol{W}} \tag{41}$$

where $\boldsymbol{\Phi}(\delta) \in \mathbb{R}^{nk}$ is the $k$-times concatenated basis function vector, $n = 3$ as $\delta \in \mathbb{R}^3$, and $\hat{\boldsymbol{W}}$ is a critic NN weight vector to be determined. Using (41), the *policy evaluation* step of the IRL Algorithm 1 can be re-written as

$$e_b = \hat{\boldsymbol{W}}^T \triangle \boldsymbol{\Phi}(\delta) - \kappa(t) \tag{42}$$

where $e_b$ is the continuous-time counterpart residual error of the TD, $\triangle \boldsymbol{\Phi}(\delta) = \boldsymbol{\Phi}(\delta(t)) - \boldsymbol{\Phi}(\delta(t + T))$, and reinforcement

$$\kappa(t) = \int_t^{t+T} \{\delta^T \boldsymbol{Q} \delta + U(\pi^p(\delta)) - U(\pi^e(\delta))\} d\tau. \tag{43}$$

Therefore, (42) implies that the problem of solving the HJI equation is converted to tuning the critic NN weights such that $e_b$ to be minimized. Now, to adjust these weights, the following objective function is employed

$$E_b = \frac{1}{2} e_b^2. \tag{44}$$

Then, the TD gradient descent algorithm[30] to minimize $e_b$ is obtained by using the chain rule

$$\dot{\hat{\boldsymbol{W}}} = -\frac{\alpha_L \triangle \boldsymbol{\Phi}(\delta)}{(1 + \triangle \boldsymbol{\Phi}(\delta)^T \triangle \boldsymbol{\Phi}(\delta))^2} e_b \tag{45}$$

where $\alpha_L > 0$ is the learning rate. The proof of convergence of critic NN weights is shown in the Theorem 3 of Modares et al..[30]

## 6 | GENERALIZED ROTATIONAL DYNAMICS OF THE PURSUER AND EVADER

The analysis in the preceding sections has shown how to derive velocity tracker for the PE game given velocity dynamics (3). In this section, we analyze the general rotational dynamics (2) that are coupled to (1), and hence (3). We first derive the desired attitude of the system by using the *Z-Y-X* Euler angle rotation matrix from [E] (earth frame) to $[B^i]$ (body frames of pursuer or evader) as shown in Figure 1, and desired force vector $\boldsymbol{f}_d^i$ in (7). Then, by the analysis developed on the desired Euler angles, we propose the desired attitude representation with unit quaternions to overcome the singularity

problem of the Euler angles. Lastly, by treating forces and/or moments as final inputs to the Newtonian system, we close the backstepping control loop to track desired force vector $\boldsymbol{f_d^i}$ in (7). Note that in this section, $i$ represents either $p$ or $e$.

Assume that the gravity $g$ is constant and the Earth is flat in the 3-dimensional space as illustrated in the Figure 1. Then, the vehicle carrier frame is aligned with the body frame $[B^i]$. Thereby the rotation matrix from $[E]$ to $[B^i]$ frames shown in Figure 1, can be given in terms of the Euler angles as

$$N(\boldsymbol{\eta^i}) = \begin{bmatrix} c\theta^i c\psi^i & c\theta^i s\psi^i & -s\theta^i \\ -c\varphi^i s\psi^i + s\varphi^i s\theta^i c\psi^i & c\varphi^i c\psi^i + s\varphi^i s\theta^i s\psi^i & s\varphi^i c\theta^i \\ s\varphi^i s\psi^i + c\varphi^i s\theta^i c\psi^i & -s\varphi^i c\psi^i + c\varphi^i s\theta^i s\psi^i & c\varphi^i c\theta i \end{bmatrix} \tag{46}$$

where $c$ and $s$ refers to cosine and sine respectively, and $\boldsymbol{\eta^i} = \begin{bmatrix} \psi^i & \theta^i & \varphi^i \end{bmatrix}^T$ is the Euler angle vector. Note that $\boldsymbol{N(\eta^i)}$ belongs to the special orthogonal group and is of rank 3, or $SO(3)$, whose determinant is equal to 1.

Assuming the direction of the thrust force to be along the nose of players' bodies or positive $x_i$-axis ($\forall i = p, e$). This enables us to write that the desired force vector is indeed in the form of $\boldsymbol{f_{B_d}^i} = \begin{bmatrix} \mu_d^i & 0 & 0 \end{bmatrix}^T$, whose components written in $[B^i]$. Using (7) and expressing the desired force $\begin{bmatrix} \mu_d^i & 0 & 0 \end{bmatrix}^T$ in $[E]$ by $\boldsymbol{f_d^i} = N^T(\boldsymbol{\eta_d^i})\boldsymbol{f_{B_d}^i}$, following relation is derived

$$\boldsymbol{f_d^i} = \begin{bmatrix} f_{x_d}^i \\ f_{y_d}^i \\ f_{z_d}^i \end{bmatrix} = N^T(\boldsymbol{\eta_d^i}) \begin{bmatrix} \mu_d^i \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} \mu_d^i(c\theta_d^i c\psi_d^i) \\ \mu_d^i(c\theta_d^i s\psi_d^i) \\ \mu_d^i(-s\theta_d^i) \end{bmatrix} \tag{47}$$

where $\boldsymbol{\eta_d^i} = \begin{bmatrix} \psi_d^i & \theta_d^i & \varphi_d^i \end{bmatrix}^T$ is the desired Euler angle vector.

Then, (47) can be solved for desired attitude angles $\theta_d$, $\psi_d$, and $\mu_d$ as

$$\theta_d^i = -\tan^{-1}\left( \frac{f_{z_d}^i}{f_{x_d}^i \cos\psi_d^i + f_{y_d}^i \sin\psi_d^i} \right), \tag{48}$$

$$\psi_d^i = \tan^{-1}\left( \frac{f_{y_d}^i}{f_{x_d}^i} \right), \tag{49}$$

$$\mu_d^i = \sqrt{f_{x_d}^{i2} + f_{y_d}^{i2} + f_{z_d}^{i2}}. \tag{50}$$

Note that $\varphi_d^i$ can be arbitrarily prescribed. However, (48)–(50) assumes that the equality conditions $f_{x_d}^i = 0, f_{y_d}^i = 0$ cannot occur simultaneously since (48) and (49) become indefinite. This singularity problem is also known as gimbal lock, which is associated with $\theta_d^i = \pi/2$.

To avoid gimbal lock, define the following unit quaternion representation

$$\boldsymbol{q^i} = \begin{bmatrix} q_0^i & q_1^i & q_2^i & q_3^i \end{bmatrix}^T = \begin{bmatrix} q_0^i & \boldsymbol{q_v^i}^T \end{bmatrix}^T \tag{51}$$

$$q_0^i = \cos\phi^i/2 \tag{52}$$

$$\boldsymbol{q_v^i} = \boldsymbol{k^i}\sin\phi^i/2 \tag{53}$$

where $\phi^i$ is the rotation about equivalent axis $\boldsymbol{k^i}$, which is subjected to constraint $\boldsymbol{q^i}^T\boldsymbol{q^i} = 1$. Moreover, the kinematics equation for unit quaternion is

$$\dot{\boldsymbol{q}}^i = \frac{1}{2}J^T(\boldsymbol{q^i})\boldsymbol{w_B^i} \tag{54}$$

where $J(\boldsymbol{q^i}) \in \mathbb{R}^{3\times4}$ satisfies the equalities $J(\boldsymbol{q^i})J^T(\boldsymbol{q^i}) = I_{3\times3}, J(\boldsymbol{q^i})\boldsymbol{q^i} = \boldsymbol{0}$, and can be expressed as

$$J(\boldsymbol{q^i}) = [-\boldsymbol{q_v^i} \quad \boldsymbol{S}(\boldsymbol{q_v^i}) + q_0^i \boldsymbol{I_{3\times 3}}]$$

$$\text{where } \boldsymbol{S}(\boldsymbol{q_v^i}) = \begin{bmatrix} 0 & q_3 & -q_2 \\ -q_3 & 0 & q_1 \\ q_2 & -q_1 & 0 \end{bmatrix}. \tag{55}$$

Then, the rotation matrix from $[B^i]$ to $[E]$ in terms of the unit quaternion (51) is given by

$$\boldsymbol{N^T}(\boldsymbol{q^i}) = \boldsymbol{I_{3\times 3}} - 2q_0^i \boldsymbol{S}(\boldsymbol{q_v^i}) + 2\boldsymbol{S^2}(\boldsymbol{q_v^i}), \tag{56}$$

which is also known as *Rodrigues* formula. The following set of equations can be obtained by substituting the rotation matrix with the argument $\boldsymbol{q_d^i}$ (56) into (47) along with selected $\varphi_d^i$

$$\begin{bmatrix} f_{xd}^i \\ f_{yd}^i \\ f_{zd}^i \end{bmatrix} = \mu_d^i \begin{bmatrix} 1 + 2(-q_{2_d}^{i\,2} - q_{3_d}^{i\,2}) \\ 2q_{0_d}^i q_{3_d}^i + q_{1_d}^i q_{2_d}^i \\ -2q_{0_d}^i q_{2_d}^i + q_{1_d}^i q_{3_d}^i \end{bmatrix}$$

$$\varphi_d = \tan^{-1}\left( \frac{2(q_{0_d}^i q_{1_d}^i + q_{2_d}^i q_{3_d}^i)}{1 - 2(q_{1_d}^{i\,2} + q_{2_d}^{i\,2})} \right). \tag{57}$$

Notice that $\boldsymbol{f_d^i} = \begin{bmatrix} f_{x_d}^i & f_{y_d}^i & f_{z_d}^i \end{bmatrix}^T$, $\varphi_d^i$ and $\mu_d^i$ are known by (7) and (48)–(50). Thence, (57) can be solved for the desired unit quaternion $\boldsymbol{q_d^i} = \begin{bmatrix} q_{0_d}^i & q_{1_d}^i & q_{2_d}^i & q_{3_d}^i \end{bmatrix}^T$ as (57) represents four equations with four unknowns, which are the elements of $\boldsymbol{q_d^i}$. Further substitute $\boldsymbol{q_d^i}$ into kinematics Equation 54 to find the desired angular velocity $\boldsymbol{w_{B_d}^i}$ such that

$$\boldsymbol{w_{B_d}^i} = 2\boldsymbol{J}(\boldsymbol{q_d^i})\dot{\boldsymbol{q}}_d^i. \tag{58}$$

*Remark* 6. For any Newtonian system (1) or (3) and (2), we know that forces and moments are coupled to each other, which implies that $\boldsymbol{\tau_B^i}$ is required to be compatible with the selected desired force $\boldsymbol{f_d^i}$ in (7).

Then applying the dynamic inversion technique, $\boldsymbol{\tau_B^i}$ is given using (58) as

$$\boldsymbol{\tau_B^i} = \boldsymbol{I_B^i}\dot{\boldsymbol{w}}_{B_d}^i - \boldsymbol{S}(\boldsymbol{w_{B_d}^i})\boldsymbol{I_B^i}\boldsymbol{w_{B_d}^i}. \tag{59}$$

Notice that we treat $\boldsymbol{\tau_B^i}$ as final input for the general rotational dynamics (2). Consequently, we will not develop further analysis by giving location of thrusters and actuators, which is a control allocation problem and out of scope of this article. Interested reader can check our work[23] to examine how to generate $\boldsymbol{\tau_B^i}$ for the quadrotors.

# 7 | IMPLEMENTATION ON DYNAMIC SYSTEM

This section reveals the simulation results of ZS PE game with different scenarios. First, we consider when both the pursuer and evader follows their game optimal velocities given in (17) and (18) respectively. Then, we show the scenario in which the pursuer tracks its game optimal velocity (17) whereas the evader adopts a sub-optimal velocity policy.

In order to model the constrained optimal velocity trajectories, (13) is evaluated for pursuer and evader. Then, the resultant integral is found as

$$U(\boldsymbol{v^{i^*}}) = \lambda^i (\nabla V^*)^T \tanh(\boldsymbol{v^{i^*}}) - 2\lambda^i \underline{\boldsymbol{R}}^i \log(\cosh(\boldsymbol{v^{i^*}})) \quad \forall i = p, e \tag{60}$$

where $\log(.)$ is the natural logarithm, $\underline{\boldsymbol{R}}^i \triangleq diag(\boldsymbol{R}^i) = [r_1^i \ r_2^i \ r_3^i]^T$, and $\boldsymbol{v^{i^*}}$ stands for the optimal velocity policy given by (17) and (18) $\forall i = p, e$.

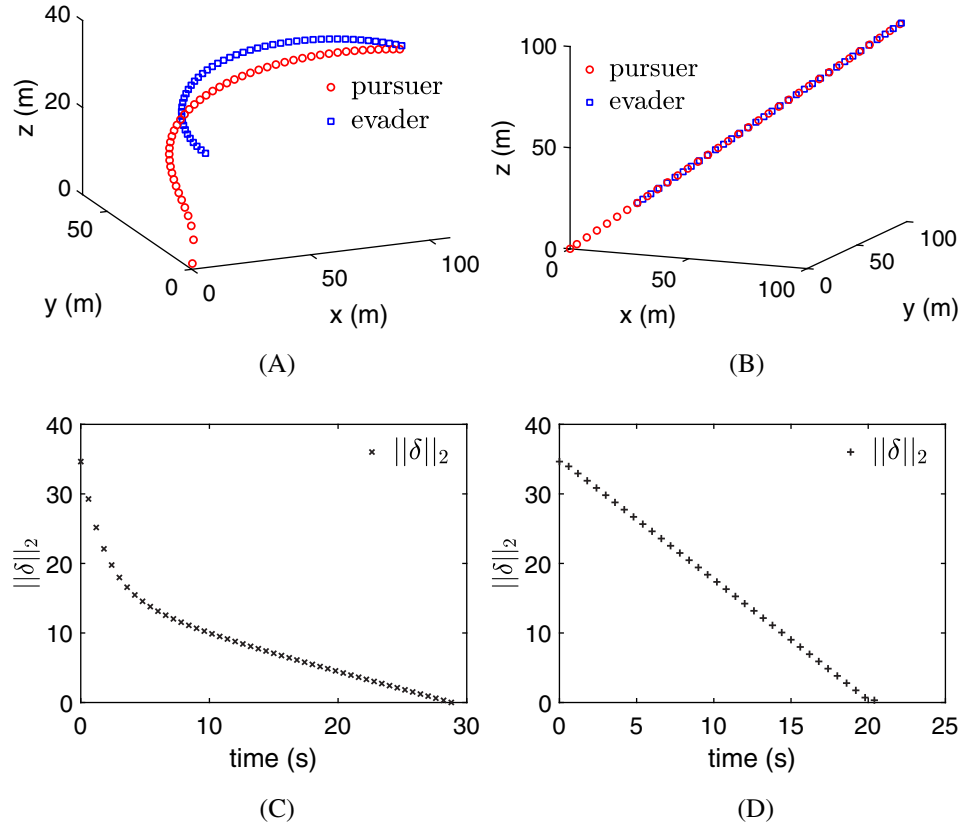**FIGURE 2** Position of the pursuer and evader: (A) $\pi^p(\delta) \triangleq$ (17), $\pi^e(\delta) \triangleq$ *suboptimal*, (B) $\pi^p(\delta) \triangleq$ (17), $\pi^e(\delta) \triangleq$ (18). $L_2$ norm of (10): (C) $\pi^p(\delta) \triangleq$ (17), $\pi^e(\delta) \triangleq$ *suboptimal*, (D) $\pi^p(\delta) \triangleq$ (17), $\pi^e(\delta) \triangleq$ (18) [Colour figure can be viewed at wileyonlinelibrary.com]

When the evader is moving with the sub-optimal velocity, we set $U(\pi^e(\delta))$ term in (11) to zero and thereby we obtain Hamilton Jacobi Bellman (HJB) equation instead of HJI (19). Notice that HJB equation in this case stands for the single player game where the pursuer is the only player. Furthermore, the existence of unique Nash equilibrium by Theorem 1 implies that the value functional (11) is convex in $\boldsymbol{\nu}^{p*}$ for $|v_j^e| \leq \lambda^e \ \forall j \in \{1,2,3\}$ given in (13), and the functional (11) is concave in $\boldsymbol{\nu}^{e*}$ for $|v_j^p| \leq \lambda^p \ \forall j \in \{1,2,3\}$. Then, (11) is separable, and solution of the HJB in terms of the optimal velocity policy for the pursuer remains the same as (17).

We conducted two simulation scenarios to validate the proposed methods in this article. We first consider ZS game with the value functional (11), and get the players track desired game optimal velocity trajectories (17), (18) by selecting ideal forces of players derived in (7) and corresponding moments (59). Then, we set $U(\pi^e(\delta))$ term in (11) to zero, and by solving the corresponding HJB equation, we played single-player game where the pursuer is the only player. Figure 2 shows the trajectories followed by the players for each of these scenarios.

In these simulations (Figures 2 and 3), parameters of the system (1), (2) are selected as $m^i = 1$ kg, $g = 9.81$ m/s$^2$, $\boldsymbol{I}_B^i = \boldsymbol{I}_{3\times3}$, where $\boldsymbol{I}_{3\times3}$ is a $3 \times 3$ identity matrix. The backstepping gain $\boldsymbol{K}^i = 5\boldsymbol{I}_{3\times3}$. In addition, the bounds (13) are $\lambda^p = 5$, $\lambda^e = 4$, and value functional parameters (11) are $\boldsymbol{Q} = 3\boldsymbol{I}_{3\times3}$, $\boldsymbol{R}^p = 0.1\boldsymbol{I}_{3\times3}$, $\boldsymbol{R}^e = 0.125\boldsymbol{I}_{3\times3}$. The position data of each player is collected through each iteration over the period $T = 0.01s$. Lastly $r^e + r^p$ (35) and shown in Figure 1 is selected as 0.25 m.

Notice that Figure 2 shows the trajectories of the players (Figure 2A,B), and corresponding $L_2$ norm of the position offset (Figure 2C,D). In addition, regarding the optimal velocity policies for the pursuer and evader, that is, when $\pi^p(\delta) \triangleq$ (17), $\pi^e(\delta) \triangleq$ (18), Figure 3 illustrate optimal velocities (17), (18), control forces (47), $L_2$ norm of velocity error (5), and Euler angles (48), (49).

Figure 4 shows the simulation results of the PE game when the velocity bounds are $\lambda^p = 10$, $\lambda^e = 9$, and other simulation parameters remain the same as in the PE game illustrated in Figures 2 and 3.
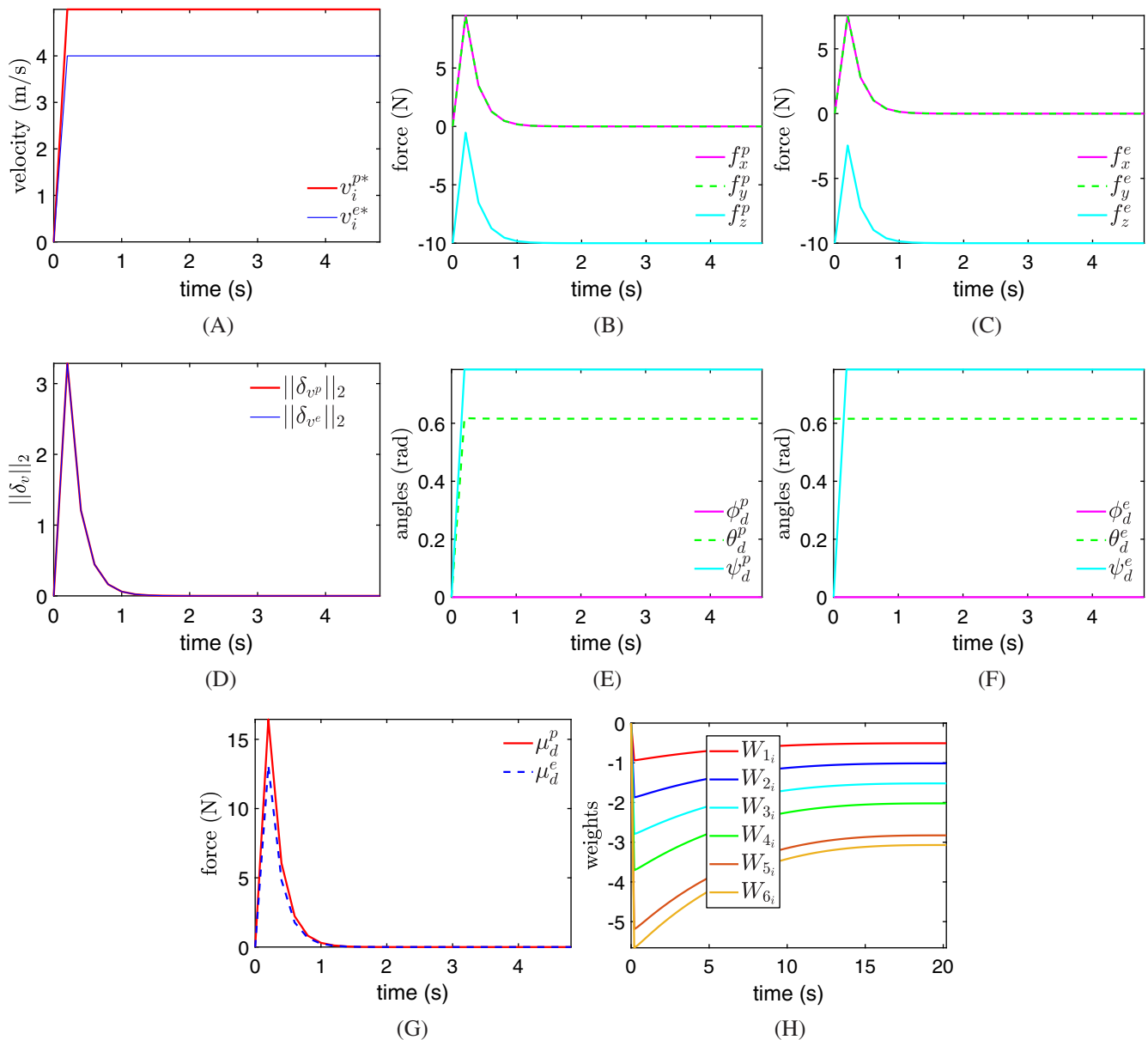
**FIGURE 3** PE game when $\lambda^p = 5$, $\lambda^e = 4$: (A) game-optimal velocity of the pursuer and evader $\forall i \in \{x, y, z\}$, (B) control force of the pursuer (47), (C) control force of the evader (47), (D) $L_2$ norm of (5) for each player, (E) Euler angles of the pursuer by (47), (F) Euler angles of the evader by (47), (G) body evaluated control force of the players (50), (H) weights $\forall i \in \{x, y, z\}$ convergence for the critic NN (45) [Colour figure can be viewed at wileyonlinelibrary.com]

# 8 | SUMMARY AND CONCLUSIONS

In this article, we worked on the game theoretic solution of pursuit-evasion (PE) intercept problem when the velocity constraints are imposed on both pursuer and evader. By solving the HJI equation corresponds to the novel non-quadratic functional, we showed that game-optimal velocity trajectories are smooth and satisfies the predetermined boundaries. Using the rigorous Lyapunov analysis, we proved that the PE game ends in a finite-time under certain conditions, which indeed implies that intercept or capture occurs in finite-time. To solve the HJI equation, the IRL method with critic NN structure is used. Consequently, we showed the simulation results of the PE game when the evader adopts both game optimal and sub-optimal velocity policy while the pursuer tracks corresponding game optimal velocity trajectory with the nonlinear backstepping tracker. Simulations showed that when the evader adopts its game optimal velocity policy, it
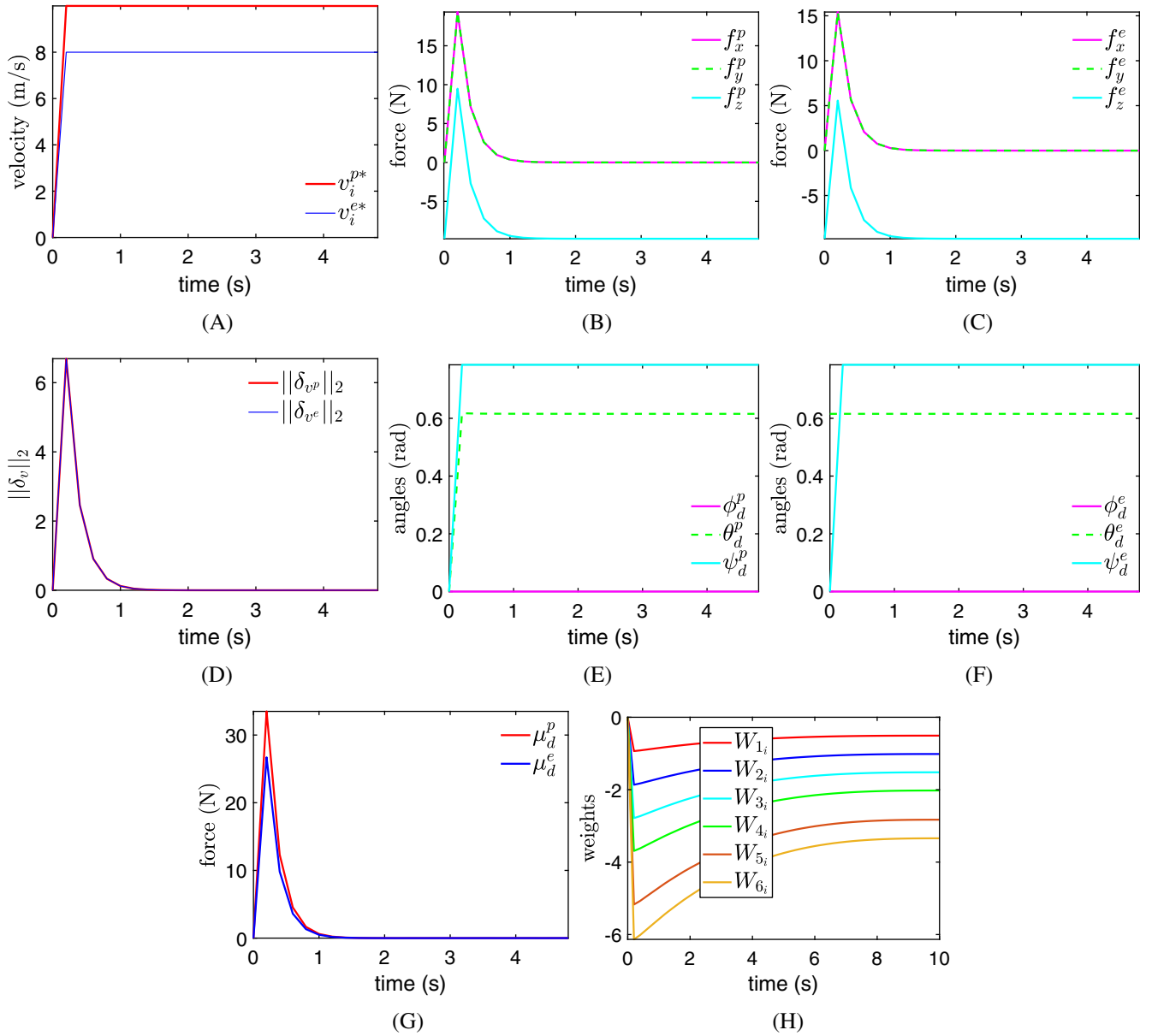
**FIGURE 4** PE game when $\lambda^p = 10$, $\lambda^e = 8$: (A) game-optimal velocity of the pursuer and evader $\forall i \in \{x, y, z\}$, (B) control force of the pursuer (47), (C) control force of the evader (47), (D) $L_2$ norm of (5) for each player, (E) Euler angles of the pursuer by (47), (F) Euler angles of the evader by (47), (G) body evaluated control force of the players (50), (H) weights $\forall i \in \{x, y, z\}$ convergence for the critic NN (45) [Colour figure can be viewed at wileyonlinelibrary.com]

takes more time to be intercepted by the pursuer compared to the scenario, in which the evader employs a sub-optimal velocity policy. Further research can be conducted to analyze robustness of the proposed methods, when the evader has mixed strategies.

## CONFLICT OF INTEREST

The authors declared that they have no conflicts of interest for this article.

## DATA AVAILABILITY STATEMENT

Data sharing is not applicable to this article as no new data were created or analyzed in this study.

## ORCID

*Yusuf Kartal* https://orcid.org/0000-0002-0295-8584
*Kamesh Subbarao* https://orcid.org/0000-0003-4295-3224

## REFERENCES

1. Turetsky V, Shinar J. Missile guidance laws based on pursuit–evasion game formulations. *Automatica*. 2003;39(4): 607-618.
2. Mylvaganam T, Sassano M, Astolfi A. A differential game approach to multi-agent collision avoidance. *IEEE Trans Autom Control*. 2017;62(8):4229-4235.
3. Marden JR, Shamma JS. Game theory and control. *Annu Rev Control Robot Autonom Syst*. 2018;1:105-134.
4. Isaacs R. *Games of Pursuit*. Santa Monica, California: Rand Corporation; 1951.
5. Isaacs R. *Differential Games: A Mathematical Theory with Applications to Warfare and Pursuit, Control and Optimization*. USA: Courier Corporation; 1999.
6. Bryson AE. *Applied Optimal Control: Optimization, Estimation and Control*. Boca Raton, FL: CRC Press; 1975.
7. Lewis FL, Vrabie D, Syrmos VL. *Optimal Control*. Hoboken, NJ: John Wiley & Sons; 2012.
8. Al-Tamimi A, Lewis FL, Abu-Khalaf M. Model-free Q-learning designs for linear discrete-time zero-sum games with application to H-infinity control. *Automatica*. 2007;43(3):473-481.
9. Basar T, Bernhard P. *H-Infinity Optimal Control and Related Minimax Design Problems: A Dynamic Game Approach*. Berlin, Germany: Springer Science & Business Media; 2008.
10. Rizvi SAA, Lin Z. Output feedback Q-learning for discrete-time linear zero-sum games with application to the H-infinity control. *Automatica*. 2018;95:213-221.
11. Hayoun SY, Weiss M, Shima T. A mixed L 2/L$\alpha$ differential game approach to pursuit-evasion guidance. *IEEE Trans Aerosp Electron Syst*. 2016;52(6):2775-2788.
12. Bhattacharya S, Başar T, Hovakimyan N. A visibility-based pursuit-evasion game with a circular obstacle. *J Optim Theory Appl*. 2016;171(3):1071-1082.
13. Li H, Liu D, Wang D. Integral reinforcement learning for linear continuous-time zero-sum games with completely unknown dynamics. *IEEE Trans Autom Sci Eng*. 2014;11(3):706-714.
14. Liu YY, Wang ZS, Shi Z. Hinf tracking control for linear discrete-time systems via reinforcement learning. *Int J Robust Nonlinear Control*. 2020;30(1):282-301.
15. Dong Y, Mingming S, Zhaowei S. Satellite proximate interception vector guidance based on differential games. *Chin J Aeronaut*. 2018;31(6):1352-1361.
16. Dong Y, Mingming S, Zhaowei S. Satellite proximate pursuit-evasion game with different thrust configurations. *Aerosp Sci Technol*. 2020;99:105715.
17. Gong H, Gong S, Li J. Pursuit-evasion game for satellites based on continuous thrust reachable domain. *IEEE Trans Aerosp Electron Syst*. 2020.
18. Jagat A, Sinclair AJ. Nonlinear control for spacecraft pursuit-evasion game using the state-dependent riccati equation method. *IEEE Trans Aerosp Electron Syst*. 2017;53(6):3032-3042.
19. Carr RW, Cobb RG, Pachter M, Pierce S. Solution of a pursuit–evasion game using a near-optimal strategy. *J Guid Control Dyn*. 2018;41(4):841-850.
20. Shaferman V, Shima T. Cooperative differential games guidance laws for imposing a relative intercept angle. *J Guid Control Dyn*. 2017;40(10):2465-2480.
21. Weintraub I, Garcia E, Pachter M. Optimal guidance strategy for the defense of a non-manoeuvrable target in 3-dimensions. *IET Control Theory Appl*. 2020;14(11):1531-1538.
22. Stevens BL, Lewis FL, Johnson EN. *Aircraft Control and Simulation: Dynamics, Controls Design, and Autonomous Systems*. Hoboken, NJ: John Wiley & Sons; 2015.
23. Kartal Y, Subbarao K, Gans NR, Dogan A, Lewis F. Distributed backstepping based control of multiple UAV formation flight subject to time delays. *IET Control Theory Appl*. 2020;14(12):1628-1638.
24. Abu-Khalaf M, Lewis FL. Nearly optimal control laws for nonlinear systems with saturating actuators using a neural network HJB approach. *Automatica*. 2005;41(5):779-791.
25. Haddad WM, Chellaboina V. *Nonlinear Dynamical Systems and Control: A Lyapunov-Based Approach*. Princeton, NJ: Princeton University Press; 2011.
26. Cannarsa P, Sinestrari C. *Semiconcave Functions, Hamilton-Jacobi Equations, and Optimal Control*. Vol 58. Berlin, Germany: Springer Science & Business Media; 2004.
27. Lopez VG, Lewis FL, Wan Y, Sanchez EN, Fan L. Solutions for multiagent pursuit-evasion games on communication graphs: finite-time capture and asymptotic behaviors. *IEEE Trans Autom Control*. 2019;65(5):1911-1923.

28. Modares H, Sistani MBN, Lewis FL. A policy iteration approach to online optimal control of continuous-time constrained-input systems. *ISA Trans.* 2013;52(5):611-621.

29. Vamvoudakis KG, Vrabie D, Lewis FL. Online adaptive algorithm for optimal control with integral reinforcement learning. *Int J Robust Nonlinear Control.* 2014;24(17):2686-2710.

30. Modares H, Lewis FL. Optimal tracking control of nonlinear partially-unknown constrained-input systems using integral reinforcement learning. *Automatica.* 2014;50(7):1780-1792.

31. Yang Y, Vamvoudakis KG, Modares H. Safe reinforcement learning for dynamical games. *Int J Robust Nonlinear Control.* 2020;30(9):3706-3726.

32. Jiang H, Zhang H, Xie X. Critic-only adaptive dynamic programming algorithms' applications to the secure control of cyber–physical systems. *ISA Trans.* 2020;104:138-144.