

# Online and Robust Intermittent Motion Planning in Dynamic and Changing Environments

Zirui Xu<sup>1</sup>, *Student Member, IEEE*, George P. Kontoudis<sup>2</sup>, *Member, IEEE*,  
Kyriakos G. Vamvoudakis<sup>3</sup>, *Senior Member, IEEE*

**Abstract**—In this paper, we propose RRT-Q<sub>∞</sub><sup>x</sup>, an online and intermittent kinodynamic motion planning framework for dynamic environments with unknown robot dynamics and unknown disturbances. We leverage RRT<sup>x</sup> for global path planning and rapid replanning to produce waypoints as a sequence of boundary value problems (BVPs). For each BVP, we formulate a finite-horizon, continuous-time zero-sum game, where the control input is the minimizer, and the worst-case disturbance is the maximizer. We propose a *robust intermittent Q-learning* controller for waypoint navigation with completely unknown system dynamics, external disturbances, and intermittent control updates. We execute a relaxed persistence of excitation technique to guarantee that the Q-learning controller converges to the optimal controller. We provide rigorous Lyapunov-based proofs to guarantee the closed-loop stability of the equilibrium point. The effectiveness of the proposed RRT-Q<sub>∞</sub><sup>x</sup> is illustrated with Monte-Carlo numerical experiments in numerous dynamic and changing environments.

## I. INTRODUCTION

Safe motion planning of mobile robots in dynamic environments is a challenging task with applications in autonomous vehicles and human-crowded robotic navigation. One difficulty of safe motion planning in dynamic environments lies in the limited ability of perception and inference, which leads to inaccuracies in multiple aspects of understanding the environment. Examples include inaccurate motion prediction for obstacles or other agents, false representation map of the environment, and imprecise localization of the ego agent. Whenever the environment model is updated, a fast replanning is needed. Moreover, optimal control is desired for autonomous systems, yet it typically requires extensive, cumbersome offline computation. In addition, precomputed optimal control suffers from the uncertain nature of the system dynamics, which compromises its optimality and safety. External disturbances are usually stochastic and unknown *a priori*. In practice, rejecting external disturbances is essential not only for optimality of the control policy, but also for maintaining safety in planning. Furthermore, mobile robots are often equipped with limited on-board computation and

communication capabilities—their energy resources are limited and require judicious distribution. In this work, our focus is on providing a safe, online, and model-free kinodynamic motion planning methodology in finite horizon without any offline computation and with intermittent communication.

Sampling-based algorithms, e.g., the Rapidly-exploring Random Tree (RRT) [1], are efficient in high-dimensional motion planning tasks. While RRT is probabilistically complete, RRT\* provides solutions that are probabilistically complete as well as asymptotically optimal in static environments by rewiring the search tree [2]. In dynamic environments, the Execution-extended RRT (ERRT) is considered as one of the first algorithms for real-time replanning [3]. While ERRT rebuilds a new tree from the current robot configuration to goal when replanning, the Dynamic RRT (DRRT) [4] only rebuilds part of the tree from the colliding nodes to goal and thus reports higher efficiency in execution time. The authors in [5] present the RRT<sup>x</sup>, an asymptotically optimal motion planning algorithm for both static and dynamic environments. Growing from goal to start, RRT<sup>x</sup> keeps the main spanning tree while replanning. During replanning, while also performing the rewiring operation as RRT\* to guarantee asymptotic optimality, RRT<sup>x</sup> improves efficiency by maintaining a limited neighborhood size and having a relaxed rewiring cascade.

For kinodynamic motion planning, optimization-based techniques are used in kinodynamic RRT\* [6] where the optimality of paths is ensured for controllable linear systems. Yet, this open-loop approach is vulnerable to disturbances, model inaccuracies, and requires offline computation. Another kinodynamic motion planning technique is the Linear Quadratic Regulator Trees (LQR-Trees) presented in [7]. This approach employs convex optimization tools to compute Lyapunov functions and regions of attraction for feedback motion planning in linearized systems. A similar approach [8] that combines LQR and RRT\*, considers time as an additional dimension of the state space. The authors in [9] present RRT-Q\*, an online, model-free kinodynamic motion planning framework that computes approximately optimal control policies for motion planning in static environments. The Stable Sparse RRT (SST) and SST\* are introduced in [10] for asymptotically near-optimal and optimal sampling-based kinodynamic motion planning, respectively. They propagate the system dynamics forward-in-time using Monte-Carlo-based random control and random propagation time without solving the boundary-value problems (BVP). Similarly to RRT, SST and SST\* provide only high-level paths without low-level control laws. A multi-query algorithm for asymptotically optimal kinodynamic mo-

<sup>1</sup>Z. Xu is with the Department of Aerospace Engineering, University of Michigan, Ann Arbor, MI 48109, USA (email: ziruiXu@umich.edu).

<sup>2</sup>G. P. Kontoudis is with the Department of Aerospace Engineering, University of Maryland, College Park, MD 20742, USA (email: kont@umd.edu).

<sup>3</sup>K. G. Vamvoudakis is with the Daniel Guggenheim School of Aerospace Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA (email: kyriakos@gatech.edu).

This work was supported, in part by NSF under grant Nos. CAREER CPS-1851588, CPS-2038589, CPS-2227185, SATC-1801611, by NASA ULI under grant number 80NSSC20M0161, and by the United States Army under Contract No. W56HZV-17-C-0095.

tion planning that sample edges instead of configurations is proposed in [11]. A survey on sampling-based motion planning algorithms appears in [12] and a benchmark in [13].

The authors in [14] combine a variant of deep reinforcement learning (RL) with RRT and propose the RL-RRT for kinodynamic motion planning. They first use deep RL to learn an obstacle avoiding policy with the system dynamics, serving as a local planner and controller. Based on the RL policy, they use supervised learning to predict the time to reach a state and guide the growth of the tree. However, this method requires extensive offline computations for policy training and still needs higher accuracy for safe navigation. In [15], a geometric version of kinodynamic motion planning is proposed. The authors combine sampling-based motion planning techniques in static environments and model-based adaptive control, to develop a robust trajectory tracking controller. The authors in [16] construct an optimal motion planning algorithm for static environments by employing RL [17] and artificial harmonic potential fields [18]. A probabilistic roadmap method of kinodynamic motion planning with known trajectories of moving obstacles is presented in [19]. This method samples new nodes in the state-time space and then maps them into the physical state space, using random control inputs of known dynamics. In [20], the authors propose an online kinodynamic motion planning framework which includes a replanning scheme to reduce replanning events in dynamic environments. Also, offline machine learning techniques drastically reduce the number of BVPs and facilitate the online implementation. However, the feasibility of this technique depends solely on the offline training of the reachability sets. The work of funnel libraries [21] employs similar convex optimization tools as in the LQR-Trees, yet incorporates the influence of bounded model uncertainties and disturbances. Online kinodynamic motion planning is addressed in unknown environments after computing offline the computationally expensive funnel libraries. The authors in [22] present a sampling-based kinodynamic motion planning framework that uses supervised learning to train a neural network controller for solving BVPs offline and apply it in real-time. The authors in [23] propose an end-to-end deep-learning-based kinodynamic motion planning, where they use a deep neural network to predict waypoints and use model predictive control to drive the system to the next waypoint. A real-time feedback motion planning and replanning method for nonlinear systems in dynamic environments with a novel graph data structure is proposed in [24]. A two-player zero-sum game (TPZSG) is formulated in [25] for the worst-case disturbance rejection in controller synthesis. Bounded rationality [26]–[28] has been combined with RRT- $Q^X$  [29] to perform multi-agent, human-like motion planning [30].

These motion planning works have reported outstanding results in various applications. Yet, almost all of them are *model-based* methods, requiring accurate knowledge of system dynamics and external disturbances for robust control synthesis. Besides, all motion planning methods are executed with constant update of the controller, leading to significant communication requirements. The constant update of control laws is not always necessary, especially when optimality can be guaranteed with an intermittent update.

Approximate dynamic programming [17] serves as a connection between adaptive control [31] and optimal control [32] by employing the principles of RL. In [33], a partially model-free algorithm is introduced to solve the optimal control problem online via policy iteration. In [34], a Q-learning technique is presented to solve the infinite-horizon optimal control problem with unknown continuous-time linear systems. The authors in [35] developed an intermediate method between policy iteration and value iteration for discrete-time systems.

Intermittent control is able to operate optimally in realistic systems with limited communication between sensors, controllers, and actuators [36]. The systems evaluate a user-defined triggering condition to determine whenever the loop should be closed. The decision is determined by the equilibrium stability of the closed-loop system. The authors in [37] propose a Q-learning approach with event-trigger conditions to solve an infinite-horizon game-theoretic problem.

The authors in [38] employ experience replay mechanisms to reduce the non-stationarities and instability of the reinforcement learning algorithm. The latter captures previous data to alleviate the update process. A similar constraint appears in adaptive control [31] that the probing signal is required to be persistently exciting for convergence of weight parameters to the ideal values. That condition is conservative and is hard to accomplish in real-world applications. In [39], [40] the authors utilize previous experiences concurrently with current data to relax the condition of persistence of excitation, which makes it easier to implement in practice. Experience replay is more recently leveraged to relax the persistence of excitation requirement in Hamiltonian-driven methods [41].

Preliminary results of this work are presented in [29], [42]. In this paper, we extend the mathematical formulation to systems with external disturbances as a two-player zero-sum game in Sections II and III and propose a robust intermittent Q-learning method in Section IV to solve it. We also provide an in-depth discussion of the algorithm along with a computational complexity analysis in Section V. In Section VI we illustrate the efficiency of the methodology in more challenging dynamic environments and present the variation of the maximum kinodynamic distances in different scenarios. Finally, in the Appendix we provide a complete Lyapunov-based stability proof by considering external disturbances.

**Contribution.** The contribution of this paper is threefold. First, we develop a *robust intermittent Q-learning* method for solving a finite-horizon game-theoretic control problem for continuous-time linear systems, without any knowledge of the system dynamics (Sections III to V). Compared to existing control techniques, our method can achieve robust intermittent control with unknown dynamics and a relaxed persistence of excitation for the finite-horizon optimal control problem. Then, we provide rigorous stability, robustness, and convergence guarantees (Theorems 2 and 3). Finally, we propose a decoupled, real-time motion planning framework, RRT- $Q_\infty^X$ , which combines robust intermittent Q-learning and a sampling-based motion planner. The proposed RRT- $Q_\infty^X$  can achieve safe, online, robust kinodynamic motion planning in unpredictable dynamic environments without offline computation/training (Table I).

**Structure.** In Section II we formulate the problem. Section III provides an intermittent optimal solution to the TPZSG. The robust intermittent Q-learning method is shown in Section IV. In Section V we discuss the algorithmic details of the proposed technique. Section VI illustrates the efficacy of the proposed motion planning method in dynamic environments, and Section VII concludes the paper.

**Notation and Nomenclature.** See the following.

$\mathbb{R}, \mathbb{N}$	real numbers, natural numbers
$\mathbb{R}^+, \mathbb{R}^{n \times m}$	positive real numbers, $n \times m$ real matrices
$\underline{\lambda}(A), \bar{\lambda}(A)$	minimum, maximum eigenvalues of matrix $A$
$\ v\ $	Euclidean norm of vector $v$
$\text{vech}(A)$	half-vectorization of matrix $A$
$\text{vec}(A)$	vectorization of matrix $A$
$\text{mat}(A)$	matrization of matrix $A$
$A[a:b, c:d]$	submatrix of row $a$ to $b$ and column $c$ to $d$ of matrix $A$
$\kappa(\cdot)$	a class $\mathcal{K}$ function
$\otimes$	Kronecker product of two matrices
$\oplus$	Minkowski sum of two sets
$\mathbf{0}_p$	vector of length $p$ with all elements 0
$\nabla_x(y)$	gradient of $y$ with respect to $x$
$M, R$	penalizing matrices
$\gamma$	disturbance rejection constant
$\rho$	admissible window for completing a BVP problem
$L, L_1$	Lipschitz constraints
$D_{\text{rob}}, D_{\text{rob}}^{\text{kin}}$	kinodynamic distance, maximum kinodynamic distance
$W_c, W_a, W_d$	critic, control actor, disturbance actor approx. weights
$\alpha_c, \alpha_a, \alpha_d$	gradient descent learning rates

## II. PROBLEM FORMULATION

In this section, we first present the formulation of a model-free two-player zero-sum game (TPZSG). Then, we incorporate the solution of the game into sampling-based path planning for dynamic environments.

### A. System Description

Consider the following continuous time linear time-invariant system describing the motion of the system,

$$\dot{x}(t) = Ax(t) + Bu(t) + Fd(t), \quad x(0) = x_0, \quad t \geq 0,$$

where  $x(t) \in \mathcal{X} \subseteq \mathbb{R}^n$  is a measurable state vector,  $u(t) \in \mathcal{U} \subseteq \mathbb{R}^m$  is the control input,  $d(t) \in \mathcal{D} \subseteq \mathbb{R}^q$  is the disturbance input,  $A \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times m}$ ,  $F \in \mathbb{R}^{n \times q}$  are the unknown plant, input, and disturbance matrices, respectively. We seek to use an intermittently updated control input to drive the system from an initial state  $x_0$  to a desired state  $x_r$  in a finite horizon  $T \in \mathbb{R}^+$ . Thus, we define a new state of differences  $\bar{x}(t) := x(t) - x_r$ , to obtain the new system,

$$\dot{\bar{x}}(t) = A\bar{x}(t) + Bu(t) + Fd(t), \quad \bar{x}(0) = x_0 - x_r, \quad t \geq 0. \quad (1)$$

To conserve computational resources and reduce communication efforts, we shall employ a sampled version of the state for intermittent learning, which is defined as,

$$\hat{\bar{x}}(t) := \begin{cases} \bar{x}(r_j), & t \in (r_j, r_{j+1}] \\ \bar{x}(t), & t = r_j \end{cases}$$

where  $\bar{x}(r_j)$  is the state of flow dynamics,  $\bar{x}(t)$  is the state of jump dynamics, and  $r_j$  is a strictly monotonically increasing

sequence of samples  $\{r_j\}_{j=1}^\infty$ . The trigger of the control signal update at  $t = r_j$  is based on an error gap described as,

$$e(t) := \hat{\bar{x}}(t) - \bar{x}(t). \quad (2)$$

We formulate the problem as a non-cooperative zero-sum game in a strategic form  $\mathcal{G} = (\mathcal{N}, \{\mathcal{S}_i\}_{i \in \mathcal{N}}, J)$ , where  $\mathcal{N} = \{P_1, P_2\}$  is the set of players,  $P_1$  is the control player selecting  $u_d$  that desires to minimize the cost and  $P_2$  is the disturbance player selecting  $d$  that aims to maximize the cost,  $\mathcal{S} = \mathcal{S}_1 \times \mathcal{S}_2 = \mathcal{U} \times \mathcal{D}$  is the set of strategies, and  $J : \mathcal{S} \rightarrow \mathbb{R}$  is the finite-horizon common cost function, which is formulated as,

$$J(\bar{x}; u_d, d; t_0, T) \quad (3)$$

$$= \phi(T) + \frac{1}{2} \int_{t_0}^T (\bar{x}^\top M \bar{x} + u_d^\top R u_d - \gamma^2 \|d\|^2) d\tau,$$

where  $u_d := \kappa(\hat{\bar{x}}(t))$  is the intermittent control input,  $\phi(T) := \frac{1}{2} \bar{x}^\top(T) P(T) \bar{x}(T)$  is the final cost with  $x(T)$  the free final state and  $P(T) \in \mathbb{R}^{n \times n} \succ 0$  a user-defined, symmetric, and positive-definite final Riccati matrix,  $M \in \mathbb{R}^{n \times n} \succeq 0$  and  $R \in \mathbb{R}^{m \times m} \succ 0$  are user-defined matrices penalizing the state and the control input respectively, and  $\gamma \in \mathbb{R}^+$  is a disturbance rejection constant,  $\gamma \geq \gamma^*$ , where  $\gamma^*$  is the smallest value that can stabilize system (1) [43]. In the finite-horizon boundary-value problem (BVP) we seek to drive the system from an initial state to a desired state in a finite time.

We seek to obtain a saddle-point equilibrium  $(u_d^*, d^*)$  such that,  $J(\bar{x}; u_d^*, d; t_0, T) \leq J(\bar{x}; u_d^*, d^*; t_0, T) \leq J(\bar{x}; u_d, d^*; t_0, T)$ , for all  $\bar{x}, u_d, d$ , provided by the optimization  $J(\bar{x}_0; u_d^*, d^*; t_0, T) = \min_{u_d} \max_d \{J(\bar{x}; u_d, d; t_0, T)\}$  subject to system (1). To solve the game given the cost function (3), we define a value function,

$$V^*(\bar{x}; t_0, T) := \quad (4)$$

$$\min_{u_c} \max_d \left\{ \phi(T) + \frac{1}{2} \int_{t_0}^T (\bar{x}^\top M \bar{x} + u_c^\top R u_c - \gamma^2 \|d\|^2) d\tau \right\}.$$

*Assumption 1:* The unknown pair  $(A, B)$  is controllable and the unknown pair  $(\sqrt{M}, A)$  is detectable.  $\square$

### B. Motion Planning Foundations

We now provide basic foundations for the motion planning case. We consider the known obstacle closed space as,  $\mathcal{X}_{\text{obs}} := \bigcup_{l=1}^{N_o} \mathcal{X}_{\text{obs},l} \subset \mathcal{X}$ , where  $N_o \in \mathbb{N}$  is the total number of obstacles. Thus, the free space is the open space  $\mathcal{X}_{\text{free}} = (\mathcal{X}_{\text{obs}})^c = \mathcal{X} \setminus \mathcal{X}_{\text{obs}}$ . In dynamic environments, the obstacle space  $\mathcal{X}_{\text{obs}}$  and the free space  $\mathcal{X}_{\text{free}}$  are functions of time. Denote the unpredictable change of the obstacle space as  $\Delta \mathcal{X}_{\text{obs}}$ , where  $\Delta \mathcal{X}_{\text{obs}}$  is unknown. For path planning, we employ RRT<sup>X</sup> that provides an optimal sub-tree which contains the planned path  $\pi(x_{0,k}, x_{r,k}; t) \in \mathbb{R}^{2(K \times n)}$ , with  $k = 1, \dots, K$ ,  $K \in \mathbb{N}$  the index of BVPs. Each BVP is described by its initial and desired states as a tuple  $(x_{0,k}, x_{r,k})$ . Since the obstacle space  $\mathcal{X}_{\text{obs}}$  evolves in time,  $\pi$  is also a function of time, and thus  $K$  also changes accordingly. RRT<sup>X</sup> constructs a graph  $G = (V, E)$ , with  $V$  the set of nodes and  $E$  the set of edges. As a slight abuse of notation, we shall refer to nodes  $v \in V$  as states  $x \in \mathcal{X}$ .

Next, let us provide connections of the game-theoretic formulation to the motion planning problem. For each BVP provided by RRT<sup>X</sup>, we aim to drive the unknown system to the desired state considering the worst-case disturbance  $d^*$  that seeks to maximize the value function (4). For the  $k$ -th BVP, define the *initial distance* as the distance from the initial state  $x_{0,k}$  to the desired state  $x_{r,k}$  as,

$$D_0(\bar{x}_{0,k}) := \|x_{0,k} - x_{r,k}\| = \|\bar{x}_{0,k}\|, \quad \forall \bar{x}_0 \in \mathbb{R}^n, \quad (5)$$

and the agent's *relative distance* to  $x_{r,k}$  as,

$$D(\bar{x}) := \|x - x_{r,k}\| = \|\bar{x}\|, \quad \forall \bar{x} \in \mathbb{R}^n. \quad (6)$$

Since the problem has a *free-final state*,  $x(T)$  will converge to a close neighborhood around  $x_{r,k}$  [32], [44]. Therefore, to reduce the navigation time, we consider that  $x_{r,k}$  is reached when the state vector enters its close neighborhood within a finite horizon. That is to say, when  $D(\bar{x}) \leq \rho D_0(\bar{x}_{0,k})$ , where  $\rho$  is the user-defined *admissible window*, the agent has reached the desired state  $x_{r,k}$ . Subsequently, the system will continue on the  $(k+1)$ -th problem.

Moreover, since the system dynamics are unknown, when RRT<sup>X</sup> generates the collision-free path  $\pi$ , it can only use straight-line paths as edges in  $E$ . However, the agent's actual trajectory is curved due to the dynamical constraints (1) as well as the optimal requirement (3). Thus, the actual trajectory is deviated from  $\pi$ , and collisions may occur when  $\pi$  is very close to the obstacles. To address this issue, we follow an obstacle augmentation strategy similarly to [9]. In particular, the algorithm constantly computes the *kinodynamic distance*,

$$D_{\text{rob}}(\bar{x}) := \frac{|\bar{x}_{0,k} \times \bar{x}|}{D_{0,k}}, \quad (7)$$

to capture the deviation of the agent's position from the corresponding straight path determined by  $(x_{0,k}, x_{r,k})$ . Then, an augmented obstacle space  $\mathcal{X}_{\text{obs}}^{\text{aug}}$  is obtained based on the *maximum kinodynamic distance*,  $D_{\text{rob}}^{\text{kin}}$ , via the Minkowski sum,

$$\mathcal{X}_{\text{obs}}^{\text{aug}} := \mathcal{X}_{\text{obs}} \oplus \mathcal{X}_{\text{kin}}, \quad (8)$$

where  $\mathcal{X}_{\text{kin}}$  is the space of a compact set bounded by a circle with center at the origin and radius  $D_{\text{rob}}^{\text{kin}}$ . Every time  $D_{\text{rob}}^{\text{kin}}$  is updated,  $\mathcal{X}_{\text{obs}}^{\text{aug}}$  will be updated by (8), and then RRT<sup>X</sup> will replan a path that is further from the actual obstacles with the newly-invalid nodes and their descendants pruned.

### III. INTERMITTENT TWO-PLAYER ZERO-SUM GAME

The Hamilton-Jacobi-Isaacs (HJI) equation [45] for the finite-horizon problem with respect to (1) and (4) yields,

$$\begin{aligned} \mathcal{H}(\bar{x}; u_c, d; \frac{\partial V^*}{\partial t}, \frac{\partial V^*}{\partial \bar{x}}) &= \frac{1}{2}(\bar{x}^\top M \bar{x} + u_c^\top R u_c - \gamma^2 \|d\|^2) \\ &+ \frac{\partial V^{*\top}}{\partial \bar{x}} (A\bar{x} + B u_c + F d) + \frac{\partial V^*}{\partial t}. \end{aligned} \quad (9)$$

An optimal value function for (1) can be defined as,

$$V^*(\bar{x}; t) = \frac{1}{2} \bar{x}^\top P(t) \bar{x}, \quad \forall \bar{x}, t \geq 0, \quad (10)$$

where  $P(t) \in \mathbb{R}^{n \times n} \succ 0$  is the solution to the following game differential Riccati equation,

$$\begin{aligned} -\dot{P}(t) &= P(t)A + A^\top P(t) + M + \gamma^{-2} P(t)F F^\top P(t) \\ &- P(t)B R^{-1} B^\top P(t), \quad t \geq 0. \end{aligned} \quad (11)$$

*Theorem 1:* Suppose that there exists a positive definite  $P(t)$ ,  $t \geq 0$  satisfying the game differential Riccati equation (11) with a terminal condition  $P(T) \succ 0$ . Then, the time-triggered state feedback optimal control takes the form of,

$$u_c^*(\bar{x}; t) = -R^{-1} B^\top P(t) \bar{x}, \quad \forall \bar{x}, t, \quad (12)$$

and the worst-case disturbance yields,

$$d^*(\bar{x}; t) = \gamma^{-2} F^\top P(t) \bar{x}, \quad \forall \bar{x}, t, \quad (13)$$

with saddle-point equilibrium value  $V^* = \bar{x}_0^\top P(0) \bar{x}_0$ .

*Proof.* The proof follows from [46, Corollary 17.1]. ■

Next, we consider an intermittent feedback controller which leads to significant reduction in communication when closing the loop. The optimal intermittent controller results in,

$$u_d^*(\hat{x}; t) := -R^{-1} B^\top P(t) \hat{x}, \quad \forall \hat{x}, t. \quad (14)$$

The main difference between the *optimal time-triggered control*  $u_c^*$  and the *optimal intermittent control*  $u_d^*$  is that, the former uses  $\bar{x}$  while the latter uses  $\hat{x}$ . Therefore, every time the control loop is closed and  $u_d^*$  is updated,  $u_d^* \leftarrow u_c^*$ .

*Remark 1:* As  $x$  approaches the desired reference state  $x_r$ , then  $\bar{x}$  approaches  $\mathbf{0}_n$ . Thus, the control laws  $u_c^*$ ,  $u_d^*$ , and the disturbance  $d^*$  also approach  $\mathbf{0}_m$  per (12), (13), (14), i.e.,  $u_c^*(T) \rightarrow \mathbf{0}_m$ ,  $u_d^*(T) \rightarrow \mathbf{0}_m$ , and  $d^*(T) \rightarrow \mathbf{0}_q$ . □

*Fact 1:* Since the system (1) is linear and both controllers (12), (14) are linear mappings of the state,  $u_c^*(\bar{x}; t) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $u_d^*(\hat{x}; t) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , the following inequality holds,

$$\begin{aligned} \|u_c^* - u_d^*\| &= \|R^{-1} B^\top P(t) (\bar{x} - \hat{x})\| \leq \|R^{-1} B^\top P(t)\| \|\bar{x} - \hat{x}\| \\ &\leq L(t) \|e\|, \end{aligned}$$

where  $L(t) \mapsto \mathbb{R}^+$  is a strictly positive-definite function [47]. □

*Lemma 1:* Given the *intermittent Hamiltonian* as,

$$\begin{aligned} \mathcal{H}(\bar{x}; u_d, d; \frac{\partial V^*}{\partial t}, \frac{\partial V^*}{\partial \bar{x}}) &= \frac{1}{2}(\bar{x}^\top M \bar{x} + u_d^\top R u_d - \gamma^2 \|d\|^2) \\ &+ \frac{\partial V^{*\top}}{\partial \bar{x}} (A\bar{x} + B u_d + F d) + \frac{\partial V^*}{\partial t}, \end{aligned} \quad (15)$$

the following inequality is satisfied.

$$\left\| \mathcal{H}(\bar{x}; u_d^*, d^*; \frac{\partial V^*}{\partial t}, \frac{\partial V^*}{\partial \bar{x}}) \right\| \leq \frac{\bar{\lambda}(R)}{2} L(t)^2 \|e\|^2. \quad (16)$$

*Proof.* Substitute the time-triggered Hamiltonian (9) into the intermittent Hamiltonian (15). Considering the optimal control (12) and the value function (10), we obtain,

$$\mathcal{H}(\bar{x}; u_d^*, d^*; \frac{\partial V^*}{\partial t}, \frac{\partial V^*}{\partial \bar{x}}) = \frac{1}{2} (u_c^* - u_d^*)^\top R (u_c^* - u_d^*). \quad (17)$$

Then, using Fact 1, inequality (16) holds. ■

*Theorem 2:* Let a positive-definite radially unbounded function  $V(\bar{x}; t) = \frac{1}{2}\bar{x}^\top P(t)\bar{x}$  with  $P$  satisfying the game differential Riccati equation (11). Then, the squared norm of the error gap follows,

$$\|e\|^2 \leq \frac{(1 - \beta^2)\lambda(M)\|\bar{x}\|^2 + \lambda(R)\|u_d^*\|^2 - \gamma^2\|d^*\|^2}{\lambda(R)L(t)^2},$$

where  $\beta \in [0, 1]$  is a user-defined bandwidth parameter, and the equilibrium point of the closed-loop system (1) is asymptotically stable as  $T \rightarrow \infty$  with the intermittent control (14) for all  $t \in (r_j, r_{j+1}]$ .

*Proof.* See Appendix A. ■

*Remark 2:* The selection of  $\beta$  depends on the available resources. By selecting  $\beta$  closer to 1, the condition in Theorem 2 is triggered more often, and the intermittent controller performs closer to the time-triggered controller. □

#### IV. ROBUST INTERMITTENT Q-LEARNING

In this section, we present a robust intermittent Q-learning approach to solve the model-free TPZSG described in Section III. This approach approximates the time-triggered optimal control policy in real-time.

Let us define the Q-function as,

$$\begin{aligned} \mathcal{Q}(\bar{x}; u_d, d; t) := & V^*(\bar{x}; t) + \mathcal{H}\left(\bar{x}; u_d, d; \frac{\partial V^*}{\partial t}, \frac{\partial V^*}{\partial \bar{x}}\right) \\ & - \mathcal{H}\left(\bar{x}; u_c^*, d^*; \frac{\partial V^*}{\partial t}, \frac{\partial V^*}{\partial \bar{x}}\right), \end{aligned} \quad (18)$$

where the time-triggered Hamiltonian associated with the optimal control and the worst-case disturbance vanishes, i.e.,  $\mathcal{H}(\bar{x}; u_c^*, d^*; \partial V^*/\partial t, \partial V^*/\partial \bar{x}) = 0$ , and  $\mathcal{Q}(\bar{x}; u_d, d; t) : \mathbb{R}^{n+m+q} \rightarrow \mathbb{R}^+$  is an action-dependent scalar value.

*Lemma 2:* Given the optimal value function  $V^*(\bar{x}; t)$  (4), the game-theoretic problem  $\mathcal{Q}^*(\bar{x}; u_d^*, d^*; t) := \min_{u_d} \max_d \mathcal{Q}(\bar{x}; u_d, d; t)$  has the following optimal value,

$$\mathcal{Q}^*(\bar{x}; u_d^*, d^*; t) = V^*(\bar{x}; t) + \frac{1}{2}(u_c^* - u_d^*)^\top R(u_c^* - u_d^*). \quad (19)$$

*Proof.* Substitute the intermittent Hamiltonian (17) into the Q-function (18). Since the time-triggered Hamiltonian associated with the optimal control and the worst-case disturbance is  $\mathcal{H}(\bar{x}; u_c^*, d^*; \partial V^*/\partial t, \partial V^*/\partial \bar{x}) = 0$ , the result follows. ■

We express the Q-function (18) in a compact quadratic form,

$$\begin{aligned} \mathcal{Q}(\bar{x}; u_d, d; t) &= \frac{1}{2}U^\top \begin{bmatrix} Q_{xx}(t) & Q_{xu_d}(t) & Q_{xd}(t) \\ Q_{u_d x}(t) & Q_{u_d u_d} & Q_{u_d d} \\ Q_{dx}(t) & Q_{du_d} & Q_{dd} \end{bmatrix} U \\ &:= \frac{1}{2}U^\top \bar{\mathcal{Q}}(t)U = \frac{1}{2}\text{vech}(\bar{\mathcal{Q}}(t))^\top (U \otimes U), \end{aligned} \quad (20)$$

where  $U := [\bar{x}^\top u_d^\top d^\top]^\top$  is the augmented state,  $Q_{xx}(t) = \dot{P}(t) + P(t) + M + P(t)A + A^\top P(t)$ ,  $Q_{xu_d}(t) = Q_{u_d x}^\top(t) = P(t)B$ ,  $Q_{xd}(t) = Q_{dx}^\top(t) = P(t)F$ ,  $Q_{u_d u_d} = R$ ,  $Q_{u_d d} = Q_{du_d} = 0$ , and  $Q_{dd} = -\gamma^2$ . Using the submatrices of  $\bar{\mathcal{Q}}(t)$  and the stationarity condition  $\partial \mathcal{Q}(\bar{x}; u_d, d; t)/\partial u_d = 0$ , we obtain a model-free formulation of the optimal intermittent control,

$$u_d^*(\hat{x}; t) = \arg \min_{u_d} \mathcal{Q}(\bar{x}; u_d, d; t) = -Q_{u_d u_d}^{-1} Q_{u_d x}(t) \hat{x}. \quad (21)$$

Similarly, by solving  $\partial \mathcal{Q}(\bar{x}; u_d, d; t)/\partial d = 0$  we formulate the worst-case disturbance as,

$$d^*(\bar{x}; t) = \arg \max_d \mathcal{Q}(\bar{x}; u_d, d; t) = -Q_{dd}^{-1} Q_{dx}(t) \bar{x}. \quad (22)$$

#### A. Actor-Critic Framework

We use a critic approximator augmented with past data to approximate the Q-function (18). Considering the compact quadratic form of the Q-function (20), let us define  $\nu(t)^\top W_c := \text{vech}(\bar{\mathcal{Q}}(t))/2 \in \mathbb{R}^{(n+m+q)(n+m+q+1)/2}$ , where  $\nu(t) \in \mathbb{R}^{((n+m+q)(n+m+q+1)/2) \times ((n+m+q)(n+m+q+1)/2)}$  is a universal basis function that depends explicitly on time,  $t \geq 0$ . Since the ideal weight parameters are unknown, we leverage adaptive control [31] for tuning the weights. Hence, the approximated Q-function takes the form of,

$$\hat{\mathcal{Q}}(\bar{x}; u_d, d; t) = \hat{W}_c^\top \nu(t) (U \otimes U). \quad (23)$$

The control actor to approximate the optimal intermittent control (21) is defined as  $W_a^\top \mu(t) := -Q_{u_d u_d}^{-1} Q_{u_d x}(t) = -R^{-1} Q_{u_d x}(t) \in \mathbb{R}^{m \times n}$ , where  $\mu(t) \in \mathbb{R}^{n \times n}$  is also a universal basis function that depends explicitly on time,  $t \geq 0$ . Thus, the approximated intermittent control policy yields,

$$\hat{u}_d(\hat{x}; t) = \hat{W}_a^\top \mu(t) \hat{x}. \quad (24)$$

Similarly, we define the disturbance actor for approximating the worst-case disturbance (22) as  $W_d^\top \xi(t) := -Q_{dd}^{-1} Q_{dx}(t) = \gamma^{-2} Q_{dx}(t) \in \mathbb{R}^{q \times n}$ , and  $\xi(t) \in \mathbb{R}^{n \times n}$  is a universal basis function that depends explicitly on time,  $t \geq 0$ . The approximated worst-case disturbance yields,

$$\hat{d}(\bar{x}; t) = \hat{W}_d^\top \xi(t) \bar{x}. \quad (25)$$

*Remark 3:* The critic and actor approximators (23), (24), (25) have no approximation errors since  $\nu(t)$ ,  $\mu(t)$ ,  $\xi(t)$  are universal spatio-temporal approximators. Therefore, since we use the whole space instead of just a compact set, the critic and actor approximators will converge to the optimum. □

Next, we leverage the integral reinforcement learning structure [17] to formulate the integral Bellman equation as,

$$\begin{aligned} V^*(\bar{x}(t); t) &= V^*(\bar{x}(t - \Delta t); t - \Delta t) \\ &\quad - \frac{1}{2} \int_{t-\Delta t}^t (\bar{x}^\top M \bar{x} + u_c^*{}^\top R u_c^* - \gamma^2 \|d^*\|^2) d\tau, \end{aligned} \quad (26)$$

$$V^*(\bar{x}(T); T) = \frac{1}{2} \bar{x}^\top(T) P(T) \bar{x}(T), \quad (27)$$

where  $\Delta t \in \mathbb{R}^+$  is a small, fixed time resolution. Moreover, using Lemma 2, (26), and (27) the following equations hold,

$$\begin{aligned} \mathcal{Q}^*(\bar{x}(t); u_d^*(t), d^*(t); t) &= \\ \mathcal{Q}^*(\bar{x}(t - \Delta t); u_d^*(t - \Delta t), d^*(t - \Delta t); t - \Delta t) & \\ - \frac{1}{2} \int_{t-\Delta t}^t (\bar{x}^\top M \bar{x} + u_d^*{}^\top R u_d^* - \gamma^2 \|d^*\|^2) d\tau, & \\ \mathcal{Q}^*(\bar{x}(T); u_d^*(T), d^*(T); T) &= \frac{1}{2} \bar{x}^\top(T) P(T) \bar{x}(T). \end{aligned}$$

## B. Learning with Past Data

The next step is to develop a learning framework to approximate  $\hat{Q}$ ,  $\hat{u}_d$ ,  $\hat{d}$  in (23), (24), (25). First, we need to guarantee that the probing signal is persistently exciting (PE).

*Definition 1:* A signal vector  $\Delta(t) : \mathbb{R}^+ \rightarrow \mathbb{R}^n$  is PE over the time interval  $[t, t + T_{\text{PE}}]$ ,  $T_{\text{PE}} \in \mathbb{R}^+$ , if there exists a strictly positive constant  $\zeta \in \mathbb{R}^+$  such that  $\zeta \mathbf{I} \leq \int_t^{t+T_{\text{PE}}} \Delta(\tau) \Delta(\tau)^\top d\tau$ , where  $\mathbf{I}$  is an identity matrix of appropriate dimensions.  $\square$

*Corollary 1:* If a signal vector  $\Delta(t)$  is PE, then it is guaranteed that the unknown parameter vectors of Q-learning will exponentially converge to the optimal values, i.e.,  $\hat{W}_c^\top \nu(t) \rightarrow (W_c^\top \nu)^*$ ,  $\hat{W}_a^\top \mu(t) \rightarrow (W_a^\top \mu)^*$ , and  $\hat{W}_d^\top \xi(t) \rightarrow (W_d^\top \xi)^*$ .

*Proof.* The proof follows from [31, Corollary 4.3.1].  $\blacksquare$

In practice, ensuring a PE probing signal is challenging. To this end, we inherit the analysis in [39], [40] to relax the PE condition. This approach employs past recorded data concurrently with current data for learning adaptation. Hence, if the recorded data is sufficiently rich, then the convergence of parameters to the optimal values can be guaranteed without enforcing the PE condition. In other words, initially we apply a PE signal and concurrently we store the past data of the probing signal in a buffer. When the buffer ensures a rich enough signal, we terminate the PE condition and continue the learning with the data collected in the buffer. Note that the data recording mechanism can be aperiodic.

The unknown elements of the intermittent Hamiltonian (15) are  $\hat{x}$  and  $V^*$ . Given  $\Delta t$  as the time resolution, we approximate  $\hat{x}$  using finite differences,

$$\dot{\hat{x}} \approx (\bar{x}(t) - \bar{x}(t - \Delta t)) / \Delta t. \quad (28)$$

According to Lemma 2 and (23), we approximate the derivatives of  $V^*$  as,

$$\frac{\partial V^*}{\partial \bar{x}} \approx \frac{\partial \hat{Q}}{\partial \bar{x}} = \hat{W}_c^\top \nu(t) \nabla_{\bar{x}} (U \otimes U), \quad (29)$$

$$\frac{\partial V^*}{\partial t} \approx \frac{\partial \hat{Q}}{\partial t} = \hat{W}_c^\top \nabla_t (\nu(t) (U \otimes U)), \quad (30)$$

where  $\nabla_t(\cdot)$  denotes the gradient with respect to  $t$ , and  $\nabla_{\bar{x}}(\cdot)$  with respect to  $\bar{x}$ . Therefore, with eqs. (29) and (30), the intermittent Hamiltonian (15) is approximated as,

$$\begin{aligned} \hat{\mathcal{H}}(U; \hat{W}_c; t) &= \frac{1}{2} \text{vech} \left( \begin{bmatrix} M & 0 & 0 \\ 0 & R & 0 \\ 0 & 0 & -\gamma^2 \end{bmatrix} \right)^\top (U \otimes U) \\ &+ \hat{W}_c^\top \left( \nabla_t (\nu(t) (U \otimes U)) + \nu(t) \nabla_{\bar{x}} (U \otimes U) \right) \dot{\hat{x}}. \end{aligned} \quad (31)$$

Using the integral Bellman equation, let us define the critic

approximation error  $e_{c_1} \in \mathbb{R}$  as,

$$\begin{aligned} e_{c_1}(t) &:= \hat{Q}(\bar{x}(t); \hat{u}_d(t), \hat{d}(t); t) \\ &- \hat{Q}(\bar{x}(t - \Delta t); \hat{u}_d(t - \Delta t), \hat{d}(t - \Delta t); t - \Delta t) \\ &+ \frac{1}{2} \int_{t-\Delta t}^t (\bar{x}^\top M \bar{x} + \hat{u}_d^\top R \hat{u}_d - \gamma^2 \|d\|^2) d\tau \\ &= \hat{W}_c^\top \left( \nu(t) U(t) \otimes U(t) - \nu(t - \Delta t) U(t - \Delta t) \otimes U(t - \Delta t) \right) \\ &+ \frac{1}{2} \int_{t-\Delta t}^t (\bar{x}^\top M \bar{x} + \hat{u}_d^\top R \hat{u}_d - \gamma^2 \|d\|^2) d\tau. \end{aligned}$$

To drive the system to the final states—penalized by  $P(T)$ —we define the final critic error  $e_{c_2} \in \mathbb{R}$  as,

$$e_{c_2}(t) := \frac{1}{2} \bar{x}(t)^\top P(T) \bar{x}(t) - \hat{W}_c^\top \nu(t) (U(t) \otimes U(t)).$$

Moreover, every time we record data, the corresponding buffer critic error  $e_{\text{buff},i} \in \mathbb{R}$  is defined as,  $e_{\text{buff},i}(t_i) := \hat{\mathcal{H}}(U(t_i); \hat{W}_c; t_i) - (\hat{u}_c(t_i) - \hat{u}_d(t_i))^\top R (\hat{u}_c(t_i) - \hat{u}_d(t_i)) / 2$ , which by using (31) yields,

$$\begin{aligned} e_{\text{buff},i}(t_i) &:= \frac{1}{2} \text{vech} \left( \begin{bmatrix} M & 0 & 0 \\ 0 & R & 0 \\ 0 & 0 & -\gamma^2 \end{bmatrix} \right)^\top (U(t_i) \otimes U(t_i)) \\ &+ \hat{W}_c^\top \nu(t_i) \nabla_{\bar{x}} (U(t_i) \otimes U(t_i)) \dot{\hat{x}} \\ &+ \hat{W}_c^\top \nabla_t (\nu(t_i) (U(t_i) \otimes U(t_i))) \\ &- \frac{1}{2} (\hat{u}_c(t_i) - \hat{u}_d(t_i))^\top R (\hat{u}_c(t_i) - \hat{u}_d(t_i)), \end{aligned}$$

where  $t_i$  is the data recording time. Next, we define the control actor approximation error  $e_a \in \mathbb{R}^m$  at  $t = r_j$  as,

$$e_a(t) := \hat{W}_a^\top \mu(t) \bar{x}(r_j) + \hat{Q}_{\text{uax}}^{-1} \hat{Q}_{\text{uax}}(t) \bar{x}(r_j),$$

and the disturbance actor approximation error  $e_d \in \mathbb{R}^q$  as,

$$e_d(t) := \hat{W}_d^\top \xi(t) \bar{x}(t) + \hat{Q}_{\text{dd}}^{-1} \hat{Q}_{\text{dd}}(t) \bar{x}(t).$$

We aim to drive the errors to zero by tuning the parameters  $\hat{W}_c^\top$ ,  $\hat{W}_a^\top$ ,  $\hat{W}_d^\top$  in (23), (24), (25). Thus, we apply gradient descent to the squared norms of the errors,

$$K_1(\hat{W}_c) = \frac{1}{2} \|e_{c_1}\|^2 + \frac{1}{2} \|e_{c_2}\|^2 + \frac{1}{2} \sum_{i=1}^{N_k} \|e_{\text{buff},i}\|^2, \quad (32)$$

$$K_2(\hat{W}_a) = \frac{1}{2} \|e_a\|^2, \quad (33)$$

$$K_3(\hat{W}_d) = \frac{1}{2} \|e_d\|^2, \quad (34)$$

where  $N_k \in \mathbb{N}$  is the number of data recording at the  $k$ -th two point zero-sum game.

*Remark 4:* Since  $K_1$ ,  $K_2$ ,  $K_3$  are the sum of squared norms, the corresponding optimization problems are convex. Thus, the gradient descent rule yields the global value. Moreover, solving the game (4) is transformed into solving convex optimization problems, and hence a real-time execution is ensured for the robust intermittent Q-learning controller.  $\square$

## C. Learning Framework

The learning framework is employed to drive the approximated cost, intermittent control, and disturbance to the optimal ones, and it consists of three tuning laws. We apply a normalized gradient descent method (similar to adaptive control techniques [31]) to (32) that result in closed-form critic approximation weights,

$$\dot{\hat{W}}_c = -\alpha_c \frac{\partial K_1}{\partial \hat{W}_c} = -\alpha_c \left( \frac{\sigma(t)e_{c_1}}{(1 + \sigma(t)^\top \sigma(t))^2} + \frac{\sigma_f(t)e_{c_2}}{(1 + \sigma_f(t)^\top \sigma_f(t))^2} + \sum_{i=1}^{N_k} \frac{\omega(t_i)e_{\text{buff},i}}{(1 + \omega(t_i)^\top \omega(t_i))^2} \right), \quad (35)$$

where  $\sigma_f(t) := \nu(t)(U(t) \otimes U(t))$ ,  $\sigma(t) := \nu(t)(U(t) \otimes U(t)) - \nu(t - \Delta t)(U(t - \Delta t) \otimes U(t - \Delta t))$ ,  $\alpha_c \in \mathbb{R}^+$  is the gradient descent learning rate that specifies the convergence rate of the critic, and  $\omega(t_i)$  is defined as,

$$\omega(t_i) := \nabla_t (\nu(t_i)(U(t_i) \otimes U(t_i))) + \nu(t_i) \nabla_{\bar{x}} (U(t_i) \otimes U(t_i)) \dot{\bar{x}},$$

where  $\sigma(t), \sigma_f(t), \omega(t_i) \in \mathbb{R}^{(n+m+q)(n+m+q+1)/2}$ . By following the critic tuning law (35), the sum of squared critic errors is regulated to zero.

For the update law of the control actor, we need to consider its intermittent behavior. More specifically, the intermittent controller will be updated only when the error gap (2) is large enough to activate the triggering condition, otherwise it will remain constant as the last updated value. That is precisely an impulsive system, hence we inherit the approaches in [48], [49] to describe the intermittent control actor behavior. We obtain a closed-form tuning law for the control actor approximation weights  $\hat{W}_a$  by applying gradient descent to (33) as,

$$\begin{cases} \dot{\hat{W}}_a = 0, & t \in (r_j, r_{j+1}) \\ \hat{W}_a^+ = \hat{W}_a - \alpha_a \frac{\partial K_2}{\partial \hat{W}_a} \\ \quad = \hat{W}_a - \alpha_a \frac{\mu(t)\bar{x}}{1 + (\mu(t)\bar{x})^\top (\mu(t)\bar{x})} e_a^\top, & t = r_j \end{cases} \quad (36)$$

where the convergence rate is determined by the gradient descent learning rate  $\alpha_a \in \mathbb{R}^+$ . The control actor tuning law (36) ensures that  $e_a$  converges to zero. Next, the disturbance actor tuning law in closed-form yields,

$$\dot{\hat{W}}_d = -\alpha_d \frac{\partial K_3}{\partial \hat{W}_d} = -\alpha_d \frac{\xi(t)\bar{x}}{1 + (\xi(t)\bar{x})^\top (\xi(t)\bar{x})} e_d^\top, \quad (37)$$

where  $\alpha_d \in \mathbb{R}^+$  is the gradient descent learning rate. Through the disturbance tuning law (37), the disturbance error  $e_d$  also converges to zero.

For the theoretical convergence analysis, we define the approximation weights errors of the critic, the control actor, and the disturbance actor as  $\tilde{W}_c := W_c - \hat{W}_c \in \mathbb{R}^{(n+m+q)(n+m+q+1)/2}$ ,  $\tilde{W}_a := W_a - \hat{W}_a \in \mathbb{R}^{n \times m}$ , and  $\tilde{W}_d := W_d - \hat{W}_d \in \mathbb{R}^{n \times q}$ . Thus, using (35), (36), (37) we express their dynamics as,

$$\begin{cases} \dot{\tilde{W}}_c = -\alpha_c \left( \frac{\sigma(t)\sigma(t)^\top}{(1 + \sigma(t)^\top \sigma(t))^2} + \Lambda \right) \tilde{W}_c, \\ \tilde{W}_a^+ = \tilde{W}_a - \alpha_a \frac{(\mu(t)\bar{x})(\mu(t)\bar{x})^\top}{1 + (\mu(t)\bar{x})^\top (\mu(t)\bar{x})} \tilde{W}_a \\ \quad - \alpha_a \frac{(\mu(t)\bar{x})\bar{x}^\top}{1 + (\mu(t)\bar{x})^\top (\mu(t)\bar{x})} \tilde{Q}_{\text{xu}} Q_{\text{u}}^{-1}, & t = r_j \end{cases} \quad t \in (r_j, r_{j+1})$$

$$\begin{aligned} \dot{\tilde{W}}_d &= -\alpha_d \frac{(\xi(t)\bar{x})(\xi(t)\bar{x})^\top}{1 + (\xi(t)\bar{x})^\top (\xi(t)\bar{x})} \tilde{W}_d \\ &\quad - \alpha_d \frac{(\xi(t)\bar{x})\bar{x}^\top}{1 + (\xi(t)\bar{x})^\top (\xi(t)\bar{x})} \tilde{Q}_{\text{xd}} Q_{\text{dd}}^{-1}, \end{aligned}$$

where  $\Lambda := \sum_{i=1}^{N_k} \frac{\omega(t_i)\omega(t_i)^\top}{(1 + \omega(t_i)^\top \omega(t_i))^2}$ ,  $\tilde{Q}_{\text{xu}} := \text{mat}(\tilde{W}_c)[1:n, n+1:n+m]$ , and  $\tilde{Q}_{\text{xd}} := \text{mat}(\tilde{W}_c)[1:n, n+m+1:n+m+q]$ .

*Remark 5:* After the data recording is rich enough to obtain  $\text{rank}([\omega(t_1), \omega(t_2), \dots, \omega(t_{N_k})]) \geq n+m+q$ , matrix  $\Lambda$  gets a full rank and  $\Lambda \succ 0$ . This is the relaxed PE condition which indicates that the recorded data of the PE probing signal is sufficient for convergence of  $\hat{W}_c$  to the optimum, and thus it is no longer needed to record new data afterwards [40]. Using the relaxed PE condition  $\tilde{Q}$  converges to the actual optimal cost as  $\tilde{W}_c$ ,  $\tilde{W}_a$  and  $\tilde{W}_d$  converge to their optimal solutions.

## D. Impulsive System

Since the controller has an intermittent behavior with discrete jumps, we consider it as an impulsive system. Therefore, the closed-loop dynamics of (1) with the intermittent controller (24) and the worst-case disturbance (25) take the form of,

$$\begin{aligned} \dot{\bar{x}} &= A\bar{x} + B\hat{u}_d + F\hat{d} \\ &= A\bar{x} + B(-R^{-1}Q_{\text{ux}}(t) - \tilde{W}_a^\top \mu(t))\hat{\bar{x}} \\ &\quad + F(\gamma^{-2}Q_{\text{dx}}(t) - \tilde{W}_d^\top \xi(t))\bar{x}. \end{aligned} \quad (38)$$

Let us define the augmented state that captures the flow dynamics of the system in the time interval  $t \in (r_j, r_{j+1}]$  as,  $\psi := [\bar{x}^\top \hat{\bar{x}}^\top \tilde{W}_c^\top \text{vec}(\tilde{W}_a)^\top \text{vec}(\tilde{W}_d)^\top]^\top \in \mathbb{R}^{2n+n+m+nq+(n+m+q)(n+m+q+1)/2}$ , with time derivative,

$$\dot{\psi} = \begin{bmatrix} \dot{\bar{x}} \\ \mathbf{0}_n \\ -\alpha_c \left( \frac{\sigma(t)\sigma(t)^\top}{(1 + \sigma(t)^\top \sigma(t))^2} + \Lambda \right) \tilde{W}_c \\ \mathbf{0}_{nm} \\ \Psi(\tilde{W}_d) \end{bmatrix}. \quad (39)$$

Then, the augmented state is defined as  $\psi^+ := [\bar{x}^+{}^\top \hat{\bar{x}}^+{}^\top \tilde{W}_c^+{}^\top \text{vec}(\tilde{W}_a^+)^\top \text{vec}(\tilde{W}_d^+)^\top]^\top$  for the jump dynamics at time  $t = r_j$ , with the time derivative yielding,

$$\psi^+ = \psi(t) + \begin{bmatrix} \mathbf{0}_n \\ \hat{\bar{x}}(t) - \bar{x}(t) \\ \mathbf{0}_{(n+m+q)(n+m+q+1)/2} \\ \Psi(\tilde{W}_a^+) \\ \mathbf{0}_{nq} \end{bmatrix}, \quad (40)$$

where  $\Psi(\tilde{W}_d) := \text{vec}(-\alpha_d \frac{(\xi(t)\bar{x})(\xi(t)\bar{x})^\top}{1 + (\xi(t)\bar{x})^\top (\xi(t)\bar{x})} \tilde{W}_d - \alpha_d \frac{(\xi(t)\bar{x})\bar{x}^\top}{1 + (\xi(t)\bar{x})^\top (\xi(t)\bar{x})} \tilde{Q}_{\text{xd}} Q_{\text{dd}}^{-1}), \Psi(\tilde{W}_a^+) := \text{vec}(-\alpha_a \frac{(\mu(t)\bar{x})(\mu(t)\bar{x})^\top}{1 + (\mu(t)\bar{x})^\top (\mu(t)\bar{x})} \tilde{W}_a - \alpha_a \frac{(\mu(t)\bar{x})\bar{x}^\top}{1 + (\mu(t)\bar{x})^\top (\mu(t)\bar{x})} \tilde{Q}_{\text{xu}} Q_{\text{u}}^{-1})$ .

By analyzing the impulsive system, the main stability theorem of the robust intermittent Q-learning is then presented.

*Theorem 3:* Consider the impulsive closed-loop system (38), with the critic approximator (23), the actor approximator (24), and the disturbance approximator (25), which are tuned by the tuning laws (35), (36), and (37) respectively. Then, given that  $T \rightarrow \infty$  the origin is a globally asymptotically stable equilibrium of the closed-loop system with state  $\psi$  for all initial conditions  $\psi(0)$ , if the following inequality holds,

$$\|e\|^2 \leq \frac{(1 - \beta^2)\lambda(M)\|\bar{x}\|^2 + \lambda(R)\|\hat{u}_d\|^2 - \gamma^2\|\hat{d}\|^2}{4\lambda(R)(L(t)^2 + L_1(t)^2)}, \quad (41)$$

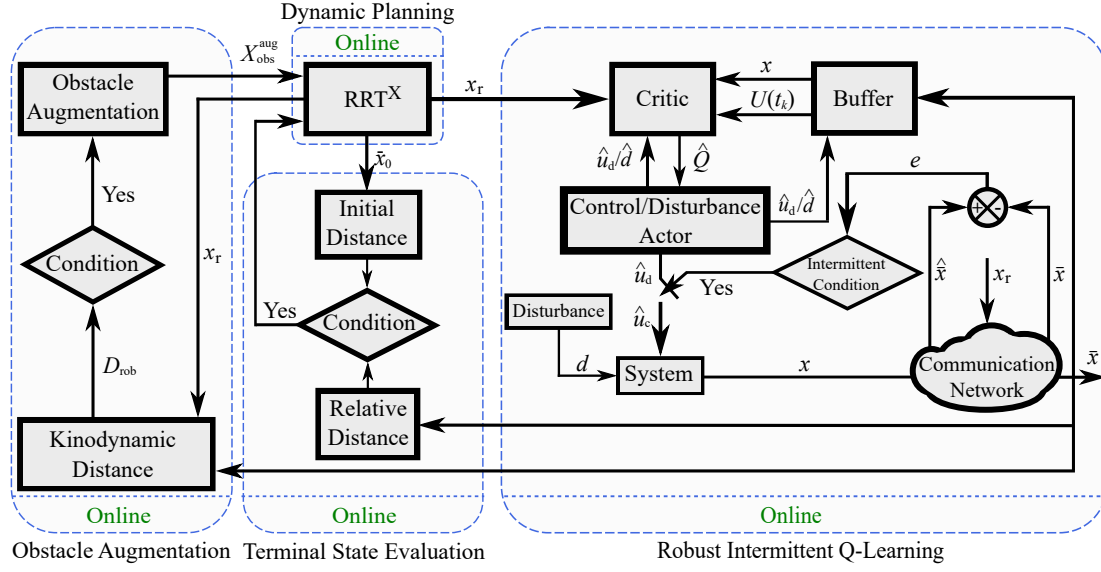


Fig. 1: **Structure of RRT-Q<sup>X</sup>**. RRT-Q<sup>X</sup> comprises four stages: i) global dynamic path planning, ii) online actor-critic framework, iii) online terminal state evaluation, and iv) online obstacle augmentation. The framework starts from global dynamic path planning and runs clockwise.

and the following inequalities of constants hold,

$$0 < \alpha_d \leq \frac{\beta^2 \lambda(M) - 2L_1(t)^2 \bar{\lambda}(R)}{\eta \bar{\lambda} \left( \frac{\xi(t) Q_{dd}^{-1}}{\|1 + (\xi(t)\bar{x})^\top (\xi(t)\bar{x})\|^2} \right)}; \quad \frac{\lambda(M)}{\bar{\lambda}(R)} > \frac{2L_1(t)^2}{\beta^2}; \quad (42)$$

$$0 < \alpha_a < \frac{8\bar{\lambda}(R) - 2}{\bar{\lambda}(R) + 2}; \quad \bar{\lambda}(R) > 0.25; \quad \alpha_c \gg \alpha_a. \quad (43)$$

*Proof.* See Appendix A. ■

## V. THE RRT-Q<sup>X</sup> FRAMEWORK

We present the RRT-Q<sup>X</sup> framework (Sections V-A and V-B) and analyze its computational complexity (Section V-C). In particular, RRT-Q<sup>X</sup> consists of four stages: i) global path planning/replanning by RRT<sup>X</sup>; ii) online control with an actor-critic framework equipped with a data buffer; iii) online obstacle augmentation; and iv) online terminal state evaluation (see Figure 1).

The main routine is presented in Algorithm 1. The gradient descent gains  $\alpha_c, \alpha_a, \alpha_d$  are designed according to Theorem 3, and the weights of approximators  $\hat{W}_c, \hat{W}_a, \hat{W}_d$  are randomly initialized (lines 1–2). In the beginning, the agent has not started moving and computes the global path using RRT<sup>X</sup> (line 20). Next, the agent performs safe navigation (lines 7–18). The agent first applies the output of the proposed robust intermittent Q-learning controller (line 7). Then, it performs the obstacle augmentation strategy to ensure a collision-free trajectory (lines 8–13). The augmented obstacle space  $\mathcal{X}_{obs}^{aug}$  is calculated by means of the Minkowski sum (8) based on the maximum kinodynamic distance  $D_{rob}^{kin}$ . Finally, the agent evaluates if it should move on to the next BVP (lines 14–18): if the agent is close enough to the desired state  $x_{r,k}$ , i.e.,  $D \leq \rho D_0$ ,  $\rho \in (0, 1)$ , then  $x_{r,k}$  will be considered as reached, and the agent will proceed to the  $(k+1)$ -th BVP. The larger the *admissible window*  $\rho$  is, the sooner the agent can proceed to the next BVP. RRT-Q<sup>X</sup> terminates once the goal state  $x_{goal}$  is considered as reached.

### Algorithm 1 RRT-Q<sup>X</sup>

**Input:**  $\mathcal{X}$  - state space;  $\mathcal{S}$  - random sample sequence;  $\mathcal{X}_{obs}$  - obstacle space;  $\Delta t$  - small time resolution;  $M, R$  - cost weight matrices;  $\gamma$  - disturbance rejection constant;  $P(T)$  - fixed final Riccati matrix;  $L, L_1$  - event-trigger constants;  $\beta$  - bandwidth parameter;  $\rho$  - admissible window;  $x_{start}$  - start state;  $x_{goal}$  - goal state;  $x(t)$  - state feedback.

**Output:**  $\hat{u}_d(t)$  - control.

```

1:  $\alpha_c, \alpha_a, \alpha_d \leftarrow \text{stability}(M, R, \gamma, L, L_1)$ ;
2:  $\hat{W}_c, \hat{W}_a, \hat{W}_d \leftarrow \text{random}$ ;  $D_{rob}^{kin} \leftarrow 0$ ;  $\mathcal{X}_{obs}^{aug} \leftarrow \mathcal{X}_{obs}$ ;  $k \leftarrow 1$ ;
3:  $V \leftarrow \{x_{goal}\}$ ;
4:  $x(0) \leftarrow x_{start}$ ;
5: while  $x(t) \neq x_{goal}$  do
6:   if agent is moving then
7:      $\hat{u}_d(t) \leftarrow \text{robustIntermittentQLearning}$ ;
8:      $D_{rob} \leftarrow \text{kinodynamicDistance}(x(t))$ ;
9:     if  $D_{rob} > D_{rob}^{kin}$  then
10:       $D_{rob}^{kin} \leftarrow D_{rob}$ ;
11:       $\mathcal{X}_{kin} \leftarrow \text{circle}(D_{rob}^{kin})$ ;
12:       $\mathcal{X}_{obs}^{aug} \leftarrow \mathcal{X}_{obs} \oplus \mathcal{X}_{kin}$ ;
13:   end if
14:    $D \leftarrow \|x(t) - x_{r,k}\|$ ;  $D_0 \leftarrow \|x_{0,k} - x_{r,k}\|$ ;
15:   if  $D \leq \rho D_0$  then
16:      $x_{0,k+1} \leftarrow x(t)$ ;
17:      $k \leftarrow k + 1$ ;
18:   end if
19: end if
20:  $G, \pi \leftarrow \text{RRT}^X(\mathcal{X}_{obs}^{aug}, \mathcal{S})$ ;
21: end while
```

#### A. Global Path Planning/Replanning by RRT<sup>X</sup>

RRT-Q<sup>X</sup> achieves global path planning and online replanning by integrating RRT<sup>X</sup> [5]. The inputs of RRT<sup>X</sup> are the agent's state space  $\mathcal{X}$  and a uniformly distributed random sample sequence  $\mathcal{S}$  of which the samples are from  $\mathcal{X}$ , and



the output is a collision-free path  $\pi$ . The waypoints of  $\pi$  will then induce a series of BVPs for the robust intermittent Q-learning control problem.

Each iteration of RRT<sup>X</sup> (Algorithm 2) begins with updating the neighborhood radius (line 4). Then, the algorithm updates the change of the obstacle space and the agent's state if any (lines 5–10). Next, the “standard” sampling-based path planner grows and refines the graph by obtaining new samples and connecting them to the existing graph if possible (lines 11–18). The rewiring operation guarantees asymptotic optimality while maintaining a limited neighborhood size (line 20).  $\epsilon$ -consistency is performed for rewiring cascade (line 21).

### B. Robust Intermittent Q-Learning

We present the robust intermittent Q-learning method (Algorithm 3). For the  $k$ -th BVP, i.e.,  $(x_{0,k}, x_{r,k})$  in the path  $\pi$ , a zero-sum game between the optimal intermittent control and the worst-case disturbance is formulated (Section III), and the proposed method performs waypoint navigation, i.e., drives the system from  $x_{0,k}$  to  $x_{r,k}$ , by solving the game. Particularly, in Algorithm 3, the critic parameters  $\hat{W}_c$  are regulated online by (35), and the critic approximates the Q-function  $\hat{Q}$  by (23) (lines 1–3). The gradient of  $\hat{W}_c$  is also influenced by `recordData` such that all weight parameters converge to the optimal values with a relaxed PE requirement (line 2). Given the approximated Q-function  $\hat{Q}$ , every time the triggering condition (41) is activated, the control actor (24) updates the approximated optimal intermittent control  $\hat{u}_d$  and the control actor parameters  $\hat{W}_a$  are tuned (36) (lines 4–8). Then, the disturbance actor parameters  $\hat{W}_d$  are tuned by (37) to approximate the worst-case disturbance  $\hat{d}$  (25) (lines 9–10). Finally, Algorithm 3 returns  $\hat{u}_d$  and applies it to the system.

With the data recording operation (Algorithm 4), the robust intermittent Q-learning method makes sure  $\hat{W}_c$ ,  $\hat{W}_a$ , and  $\hat{W}_d$  converge to their optimal values without needing the PE requirement (Corollary 1). In particular, this operation records concurrent data in a buffer and updates the gradient of  $\hat{W}_c$  until  $\text{rank}([\omega(t_1), \omega(t_2), \dots, \omega(t_{N_k})]) \geq n + m + q$ , where  $N_k$  is the number of already recorded data points. This operation can be executed either periodically or aperiodically [39], [40].

### C. Computational Complexity

The computational complexity of RRT-Q <sub>$\infty$</sub> <sup>X</sup> is determined by three components:  $\mathcal{C}_1$  for RRT<sup>X</sup>,  $\mathcal{C}_2$  for robust intermittent Q-learning, and  $\mathcal{C}_3$  for obstacle augmentation. In particular,  $\mathcal{C}_1$  for RRT<sup>X</sup> is different for static and dynamic environments. For a static environment,  $\mathcal{C}_1 = \Theta(|V| \log |V|)$  is required to build a graph  $G = (V, E)$ . When the environment changes, RRT<sup>X</sup> will first find  $V_{\text{obs}}$ , the set of all nodes whose trajectories cross the variation of the obstacle space  $\Delta \mathcal{X}_{\text{obs}}$ . Then, the expected running time will be  $\mathcal{C}_1 = \mathcal{O}(|\mathcal{D}(V_{\text{obs}})| \log |V|)$ , where  $\mathcal{D}(V_{\text{obs}})$  contains all descendants of all nodes in  $V_{\text{obs}}$  [5].

As for  $\mathcal{C}_2$  for robust intermittent Q-learning, three operations are constantly executed during the flow dynamics: approximation of Q-function, approximation of disturbance, and error gap checking. For the approximation of Q-function by (23) and (35), the complexity has a quadratic growth of  $\mathcal{O}((n+m+q)^2)$

---

#### Algorithm 2 RRT<sup>X</sup>

---

```

1:  $V \leftarrow \{x_{\text{goal}}\};$ 
2:  $x(0) \leftarrow x_{\text{start}};$ 
3: while  $x(t) \neq x_{\text{goal}}$  do
4:    $r \leftarrow \text{shrinkingBallRadius}(|V|);$ 
5:   if  $\Delta \mathcal{X}_{\text{obs}} \neq \emptyset$  then
6:      $\text{updateObstacles}(\mathcal{X}_{\text{obs}});$ 
7:   end if
8:   if agent is moving then
9:      $x(t) \leftarrow \text{measureState}();$ 
10:  end if
11:   $x_{\text{new}} \leftarrow \text{randomNode}(\mathcal{S});$ 
12:   $x_{\text{nearest}} \leftarrow \text{nearest}(x_{\text{new}});$ 
13:  if  $d(x_{\text{new}}, x_{\text{nearest}}) > \delta$  then
14:     $x_{\text{new}} \leftarrow \text{saturate}(x_{\text{new}}, x_{\text{nearest}});$ 
15:  end if
16:  if  $x_{\text{new}} \notin \mathcal{X}_{\text{obs}}$  then
17:     $\text{extend}(x_{\text{new}}, r);$ 
18:  end if
19:  if  $x_{\text{new}} \in V$  then
20:     $\text{rewireNeighbors}(x_{\text{new}});$ 
21:     $\text{reduceInconsistency}();$ 
22:  end if
23: end while
```

---



---

#### Algorithm 3 robustIntermittentQLearning

---

```

1:  $\hat{W}_c \leftarrow \text{critic}(\bar{x}, \hat{u}_d, \alpha_c, \hat{W}_c);$ 
2:  $\hat{W}_c \leftarrow \text{recordData}(\bar{x}, \bar{x}, \hat{u}_d, \hat{d}, N_k);$ 
3:  $\hat{Q} \leftarrow \text{approximateQ}(\bar{x}, \hat{u}_d, \hat{W}_c);$ 
4:  $\|e\|^2 \leftarrow \text{errorGap}(\bar{x}, \hat{x});$ 
5: if  $\|e\|^2 \geq \text{errorThreshold}(\beta, \bar{x}, \hat{u}_d, \hat{d})$  then
6:    $\hat{W}_a^+ \leftarrow \text{controlActor}(\bar{x}, \alpha_a, \hat{W}_a, \hat{Q});$ 
7:    $\hat{u}_d \leftarrow \text{approximateControl}(\bar{x}, \hat{W}_a);$ 
8: end if
9:  $\hat{W}_d \leftarrow \text{disturbanceActor}(\bar{x}, \alpha_d, \hat{W}_d, \hat{Q});$ 
10:  $\hat{d} \leftarrow \text{approximateDisturbance}(\bar{x}, \hat{W}_d);$ 
11: return  $\hat{u}_d;$ 
```

---

where  $n + m + q$  is the size of the augmented state. Similarly, the approximation of disturbance by (25) and (37) requires  $\mathcal{O}(q^2)$ . When checking the error gap by (41), the complexity is  $\mathcal{O}(n^2 + m^2 + q^2)$ . Thus, the computational complexity during the flow dynamics is  $\mathcal{O}((n+m+q)^2)$ . For the jump dynamics, the additional computation is required for updating  $\hat{W}_a$  and  $\hat{u}_d$  by (36) and (24), which takes  $\mathcal{O}(n^2 + m^2)$ . Therefore,  $\mathcal{C}_2 = \mathcal{O}((n+m+q)^2)$ .

The obstacle augmentation complexity  $\mathcal{C}_3$  depends on the Minkowski sum operation to the whole obstacle space and the following replanning. Let us assume all obstacles are convex polygons, the regular polygon that approximates the circular state space  $\mathcal{X}_{\text{kin}}$  has  $p_{\text{kin}}$  vertices, and each of the  $N_l$  obstacles in  $\mathcal{X}_{\text{obs}}$  has  $p_l$  vertices. According to [50, Theorem 13.11], the Minkowski sum of two convex polygons with  $p_1$  and  $p_2$  vertices in a plane can be computed in  $\mathcal{O}(p_1 + p_2)$  time. Then, the Minkowski sum takes  $\mathcal{O}(N_o p_{\text{kin}} + \sum_{l=1}^{N_o} p_l)$ . The following

---

**Algorithm 4** recordData

---

```

1: if rank( $[\omega(t_1), \omega(t_2), \dots, \omega(t_{N_k})]$ )  $< n + m + q$  then
2:    $\omega(t_{N_k+1}) \leftarrow \nabla_t(\nu(t)(U(t) \otimes U(t))) +$ 
    $\nu(t)\nabla_{\bar{x}}(U(t) \otimes U(t))\dot{\bar{x}};$ 
3:    $e_{\text{buff}, N_k+1} \leftarrow \text{buffErr}(\dot{\bar{x}}, \hat{u}_d, \hat{u}_c, \hat{d}, \hat{W}_c, \omega(t_{N_k+1}));$ 
4:    $\hat{W}_c \leftarrow \text{criticRelaxedPE}(\alpha_c, \hat{W}_c, e_{\text{buff}, N_k+1});$ 
5:    $N_k \leftarrow N_k + 1;$ 
6:   return  $\hat{W}_c;$ 
7: end if

```

---

replanning considering the enlargement of the whole obstacle space requires  $\mathcal{O}(|\mathcal{D}(V'_{\text{obs}})|\log|V|)$ , where  $V'_{\text{obs}}$  is the set of all the nodes whose trajectories cross the enlarged obstacle space, and  $\mathcal{D}(V'_{\text{obs}})$  contains all descendants of all nodes in  $V'_{\text{obs}}$ . Therefore,  $\mathcal{C}_3 = \mathcal{O}(N_o p_{\text{kin}} + \sum_{l=1}^{N_o} p_l + |\mathcal{D}(V'_{\text{obs}})|\log|V|)$ .

## VI. SIMULATIONS AND RESULTS

In this section, we first evaluate the robust intermittent Q-learning method in a single two-player zero-sum game (TPZSG) (Section VI-A). We next evaluate RRT-Q $^\infty$  and show its effectiveness and safety through 30 Monte-Carlo trials of the navigation task in a simulated dynamic environment with appearing/disappearing obstacles (Section VI-B). Then, we present the dependence of the kinodynamic distance on the orientation of BVP and event-trigger bandwidth (Section VI-C). Finally, we qualitatively compare RRT-Q $^\infty$  with multiple other kinodynamic motion planning works (Section VI-D).

### A. Solving a Two-Player Zero-Sum Game

We present how the robust intermittent Q-learning technique solves a single TPZSG. Consider the following linear time-invariant system that is unknown,

$$\begin{bmatrix} \dot{x}_1 \\ \dot{y}_1 \\ \dot{x}_2 \\ \dot{y}_2 \end{bmatrix} = A \begin{bmatrix} x_1 \\ y_1 \\ x_2 \\ y_2 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0.025 & 0 \\ 0 & 0.025 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} d, \quad (44)$$

$$A = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -0.5 & 0 & -1.125 & 0 \\ 0 & -0.5 & 0 & -1.125 \end{bmatrix},$$

where  $x_1, y_1$  denote the translations,  $\dot{x}_1, \dot{y}_1$  the velocities,  $\ddot{x}_1, \ddot{y}_1$  the accelerations,  $u_1, u_2$  the control inputs, and  $d$  denotes the disturbance.

We seek to drive the system from  $x_0 = [0, 3, 1, -1]$  to  $x_r = [5, 1, 0, 0]$  for  $T = 15$  s with the event-trigger bandwidth  $\beta = 0.1$  and no disturbance, i.e.,  $d = 0$ . The performance penalizing matrices are  $M = 5I_4$ ,  $R = I_2$ , the disturbance rejection constant is  $\gamma = 0.1$ , and the final Riccati matrix is  $P(T) = 0.3I_4$ . According to Theorem 3, we select the tuning gains as  $\alpha_c = 90$ ,  $\alpha_a = 0.8$ ,  $\alpha_d = 3$ , and the Lipschitz constants are  $L = 10$  and  $L_1 = 0.9\beta(\underline{\lambda}(M)/(2\bar{\lambda}(R)))^{1/2} = 0.14$ .

The evolution of states is shown in Figure 2(a) and the evolution of controls in Figure 2(b). The squared norm error (2) and its threshold (41) are presented by the blue and red lines in Figure 2(c). Every time the error (2) is close to the calculated threshold (41), the intermittent control is activated

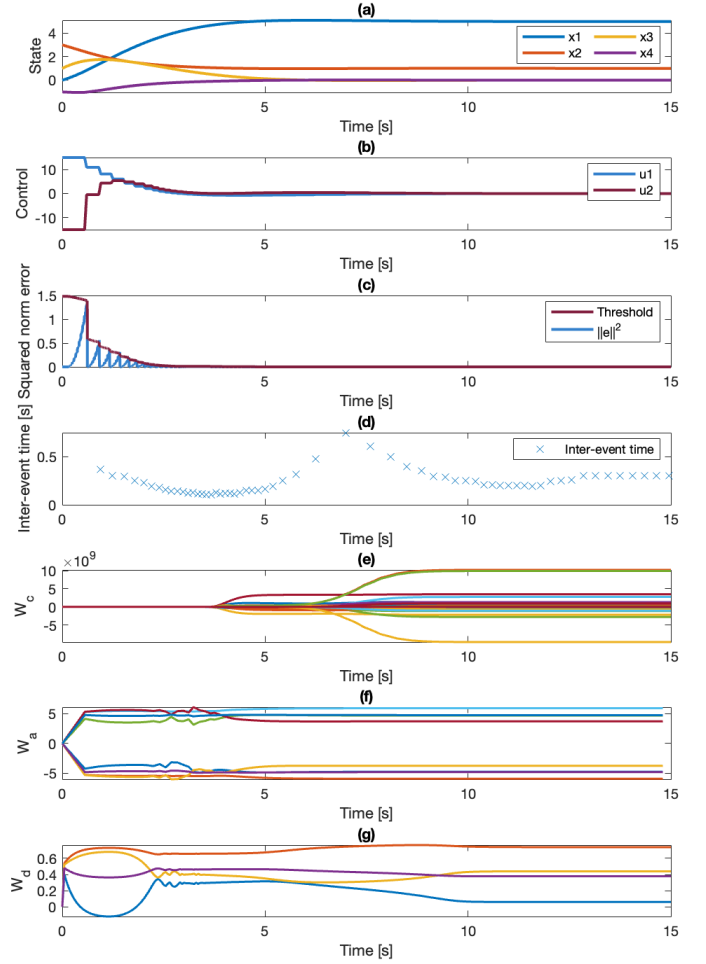


Fig. 2: **Solving a single TPZSG.** Limiting the squared norm error with a calculated threshold, the intermittent control policy guarantees the stability of the equilibrium point and conserves computational and communication efforts. (a)–(d) contain the evolution of the states, of the controls, of the error gap with the threshold, and the inter-event times of control updates for solving the game, respectively. (e)–(g) show the convergence of the approximation weights for the critic, control actor, and disturbance actor, respectively.

and the error is reset to zero. We set the upper bar of the error as 95% of the threshold, and the intermittent control is activated 61 times. The inter-event time, i.e., the time interval of two consecutive control updates, is shown in Figure 2(d). We find that the minimum inter-event time is 0.11 s and the average inter-event time is 0.24 s, which means the Zeno behavior [51] is not observed. Finally, the evolution of the approximation weights, i.e.,  $\hat{W}_c$ ,  $\hat{W}_a$ , and  $\hat{W}_d$ , is shown in Figure 2(e)–(g), where the convergence of all approximation weights is observed.

### B. Motion Planning in a Dynamic Environment

In this section, we present the kinodynamic motion planning of an autonomous rover that follows (44) in a random forest obstacle environment. We set the finite horizon to  $T = 10$  s, the event-trigger bandwidth  $\beta = 0.9$ , the performance penalizing matrices  $M = 10I_4$ ,  $R = 0.525I_2$ , the disturbance rejection constant  $\gamma = 15$ , the final Riccati matrix  $P(T) = 0.5I_4$ , and the gradient descent gains  $\alpha_c = 90$ ,  $\alpha_a = 0.8$ ,  $\alpha_d = 3$  per Theorem 3. The Lipschitz constants are  $L = 10$ ,

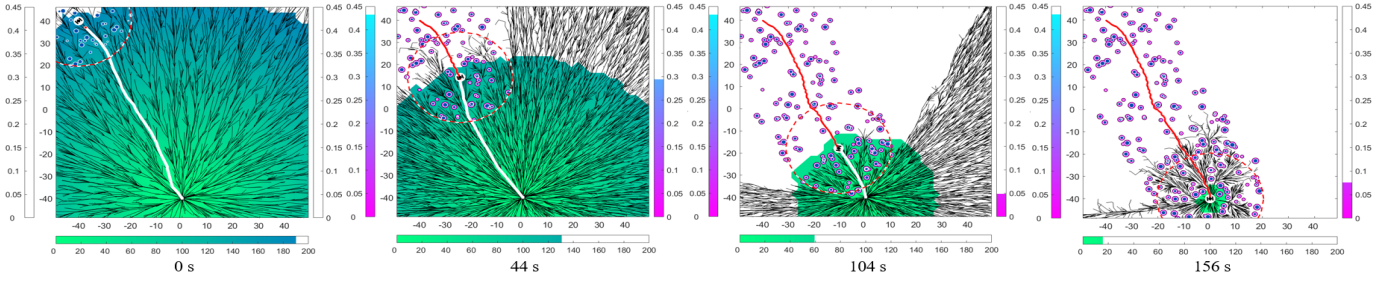


Fig. 3: **Collision-free navigation of an autonomous rover in a forest environment using RRT-Q<sup>X</sup><sub>∞</sub>**. The rover has the start state  $x_{\text{start}} = [-40, 40, 0, 0]$  and the goal state  $x_{\text{goal}} = [0, -40, 0, 0]$ . It can find the obstacles omnidirectionally within a radius of 20 represented by a red dashed circle. Its traversed path is shown by a red solid line and the RRT<sup>X</sup> path by a white line. The black solid lines represent the search-tree of RRT<sup>X</sup>. The actual obstacles (trees) are denoted by blue polygons and the corresponding augmentation in magenta. The maximum kinodynamic distance  $D^{\text{kin}}_{\text{rob}}$ , the current kinodynamic distance  $D_{\text{rob}}$ , and the *cost-to-go* of RRT<sup>X</sup> are shown in the left, right, and bottom colorbars, respectively.

$L_1 = 0.9\beta(\lambda(M)/(2\bar{\lambda}(R)))^{1/2} = 2.5$ . The initial values of  $\hat{W}_c$ ,  $\hat{W}_a$ ,  $\hat{W}_d$  are randomly generated separately for each BVP. The admissible window for each game is  $\rho = 0.2$ , i.e., when  $D(\bar{x}) \leq 0.2D_0(\bar{x}_{0,k})$ , the rover will proceed to the next BVP. In this scenario, we set the external disturbance as  $d(t) = 0.3 \sin t$ .

*Remark 6:* The larger  $\alpha_c$ ,  $\alpha_a$ , and  $\alpha_d$  are (as long as Theorem 3 is satisfied), the faster the proposed robust intermittent Q-learning controller will converge to the optimal controller (Section IV-C). Also, the larger the admissible window  $\rho$  is, the sooner the robot can proceed to the next BVP (Section V). Finally, the larger  $\beta$  is, the more frequently the control is updated (Theorem 3), and the smaller the maximum kinodynamic distance will be.  $\square$

The autonomous rover has access to state feedback and a disturbance observer. In addition, its perception range is limited, so only the obstacles within the perception range can be detected. Therefore, the obstacle space keeps enlarging while the rover navigates in the environment and perceives the newly “appearing” obstacles. RRT-Q<sup>X</sup><sub>∞</sub> measures the kinodynamic distance  $D_{\text{rob}}$  and computes the augmented obstacle space  $\mathcal{X}^{\text{aug}}_{\text{obs}}$  at every time step. Whenever the current path intersects with  $\mathcal{X}^{\text{aug}}_{\text{obs}}$ , replanning is performed to compute a new collision-free and optimal path.

We implemented RRT-Q<sup>X</sup><sub>∞</sub> partly based on the software package from [52]. The effectiveness of RRT-Q<sup>X</sup><sub>∞</sub> is shown in a video demonstration<sup>1</sup> that contains four different navigation scenarios with random forests of obstacles. Motion snapshots of one trial are depicted in Figure 3. We conduct 30 Monte-Carlo trials with different random forest obstacle configurations and different start and goal states to remove the effect of random obstacle-free navigation. The distance between the rover and the nearest obstacle in 30 replications is presented in Figure 4. The results reveal that the rover can safely navigate with no collision through the unknown random forest environment with unpredictably appearing obstacles in all trials. Moreover, the proposed controller is robust to disturbances unlike [29]. Compared to [9], [25], the proposed controller operates in open loop most of the time and closes the loop only when the error gap threshold is achieved. Hence, the controller as well as the actuator inputs are not constantly

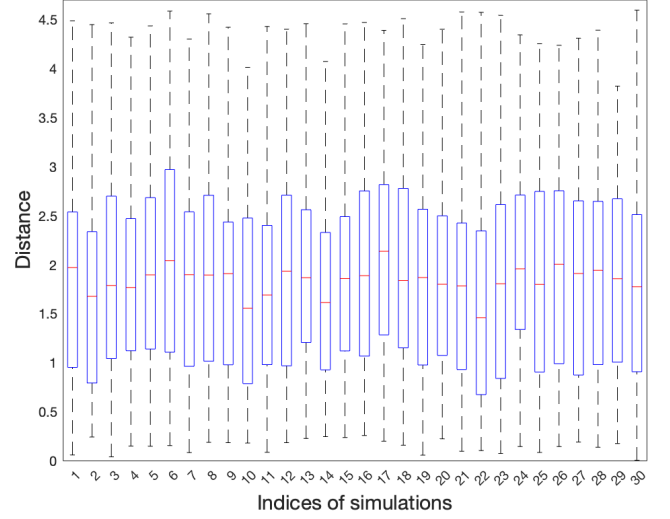


Fig. 4: **Distance to the nearest obstacle.** The distance between the robot and the nearest obstacle is positive in all 30 Monte-Carlo trials with the same navigation scenario, which shows the navigation is always collision-free.

updated. As a result, we transition from classical periodic learning to aperiodic learning that effectively balances the system performance over resources usage.

### C. Maximum Kinodynamic Distance

The value of the maximum kinodynamic distance,  $D^{\text{kin}}_{\text{rob}}$ , of a single BVP depends on multiple parameters, including the initial velocity, orientation of BVP  $\theta$  (i.e., the bearing angle of  $x_r - x_0$ ), event-trigger bandwidth  $\beta$ , as well as disturbances. We investigate how  $\theta$  and  $\beta$  are correlated with  $D^{\text{kin}}_{\text{rob}}$  with zero initial velocity and no disturbances. We implement 2,880 simulations of the rover (44) traversing from  $[0, 0]$  to  $[10 \cos \theta, 10 \sin \theta]$ . The orientation of  $\theta$  ranges from  $0^\circ$  to  $355^\circ$  with 72 values, and the event-trigger bandwidth  $\beta$  spans from 0 to 0.975 with 40 values. The other parameters are identical to the setting in Section VI-B. The colormap of  $D^{\text{kin}}_{\text{rob}}$  is presented in Figure 5. We observe that: i)  $D^{\text{kin}}_{\text{rob}}$  changes periodically with respect to  $\theta$ ; and ii)  $D^{\text{kin}}_{\text{rob}}$  decreases for  $\beta$  close to 1. Particularly, as  $\beta$  increases, the control loop is more frequently closed, which indicates more frequent regulation to the agent’s motion and subsequently smaller  $D^{\text{kin}}_{\text{rob}}$ .

<sup>1</sup>[https://youtu.be/iS\\_PzDmlpf5](https://youtu.be/iS_PzDmlpf5).

TABLE I: Kinodynamic Motion Planning Comparison

	SST & SST* [10]	Funnel Libraries [21]	Allen, Pavone [20]	RL-RRT [14]	RRT-Q* [9]	RRT-Q <sub>∞</sub> <sup>X</sup>
<b>Optimality</b>	-	Energy	Appr. Time	No	Appr. Energy	Appr. Energy
<b>Scalability</b>	-	Yes (offline computation)	Yes (training)	Yes (training)	Yes	Yes
<b>Feedback</b>	-	Closed-loop	Open-loop	Open-loop	Closed-loop	Intermittent
<b>Dynamics</b>	Model-based	Model-based (uncertain)	Model-based	Model-based	Model-free	Model-free
<b>Disturbance</b>	-	Bounded	No	No	No	Bounded
<b>Environment</b>	Static	Dynamic	Dynamic	Static	Static	Dynamic
<b>PE</b>	-	-	-	-	Strong	Relaxed

#### D. Qualitative Comparison

In Table I, we present a qualitative comparison of the proposed framework RRT-Q<sub>∞</sub><sup>X</sup> with several other works on kinodynamic motion planning. We consider seven aspects: i) optimality in terms of control; ii) scalability in terms of real-time execution; iii) feedback loop of the controller; iv) requirement for system dynamics knowledge; v) type of bearable disturbance, vi) type of environment; and vii) probing signal requirement for PE.

Although SST and SST\* [10] are near-optimal and optimal kinodynamic motion planners, the control-related specifications are not applicable due to lacking control synthesis. While the energy optimality of RRT-Q\* [9] and RRT-Q<sub>∞</sub><sup>X</sup> is approximate due to the learning process with unknown system dynamics, the time optimality of Allen and Pavone [20] arises from the prediction error of the support vector machine. Real-time computation of the control input can be achieved by all methods except of SST and SST\* [10]. Funnel libraries [21], Allen and Pavone [20], and RL-RRT [14] require extensive offline computations for training, but RRT-Q\* and RRT-Q<sub>∞</sub><sup>X</sup> require no training *a priori*. While RRT-Q\* can be used only in environments with static obstacles, RRT-Q<sub>∞</sub><sup>X</sup> can be applied in environments with appearing/disappearing obstacles. A strong PE of the probing signal is required in RRT-Q\* for parameters to converge to the ideal values while a relaxed PE condition is needed for RRT-Q<sub>∞</sub><sup>X</sup>.

#### VII. CONCLUSION

We propose RRT-Q<sub>∞</sub><sup>X</sup>, an online kinodynamic motion planning framework for unpredictable dynamic environments with unknown robot dynamics. RRT-Q<sub>∞</sub><sup>X</sup> reduces computation and communication resources by intermittently updating the controller instead of continuously and is a fully online framework with no offline computation. In particular, we propose a robust intermittent Q-learning method for low-level model-free control. The results of kinodynamic motion planning in simulated random forest environments indicate that the proposed framework enables robust collision-free navigation for completely unknown systems in unknown dynamic environments with external disturbances.

Ongoing work is focusing on: i) adaptive obstacle augmentation strategies, ii) kinodynamic motion planning in human-crowded environments, and iii) reduced-order model-free control for large-scale unknown systems.

#### REFERENCES

[1] J. J. Kuffner and S. M. LaValle, “RRT-Connect: An efficient approach to single-query path planning,” in *IEEE International Conference on Robotics and Automation*, vol. 2, 2000, pp. 995–1001.

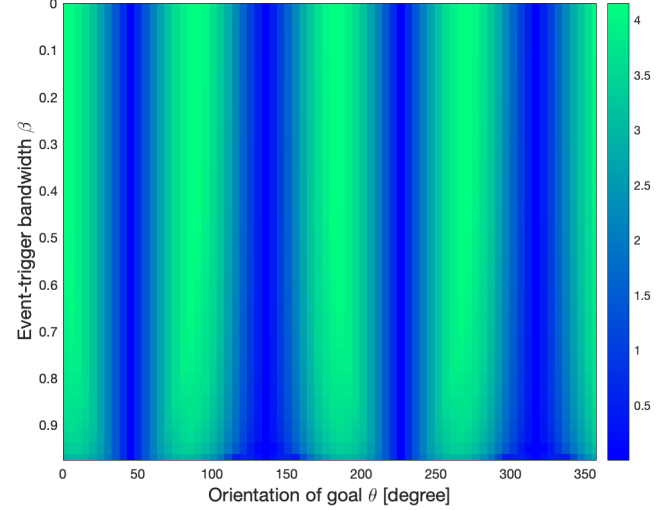


Fig. 5: **Evolution of maximum kinodynamic distance.** In 2, 880 problems over 72 values of BVP orientation  $\theta$  and 40 event-trigger bandwidth values of  $\beta$ ,  $D_{\text{rob}}^{\text{kin}}$  evolves periodically with  $\theta$  and decreases as  $\beta$  approaches 1.

[2] S. Karaman and E. Frazzoli, “Sampling-based algorithms for optimal motion planning,” *The International Journal of Robotics Research*, vol. 30, no. 7, pp. 846–894, 2011.

[3] J. Bruce and M. M. Veloso, “Real-time randomized path planning for robot navigation,” in *Robot Soccer World Cup*. Springer, 2002, pp. 288–295.

[4] D. Ferguson, N. Kalra, and A. Stentz, “Replanning with RRTs,” in *IEEE Internat. Conference on Robotics and Automation*, 2006, pp. 1243–1248.

[5] M. Otte and E. Frazzoli, “RRT<sup>X</sup>: Asymptotically optimal single-query sampling-based motion planning with quick replanning,” *The International Journal of Robotics Research*, vol. 35, no. 7, pp. 797–822, 2016.

[6] D. J. Webb and J. van den Berg, “Kinodynamic RRT\*: Asymptotically optimal motion planning for robots with linear dynamics,” in *IEEE Inter. Conference on Robotics and Automation*, 2013, pp. 5054–5061.

[7] R. Tedrake, I. R. Manchester, M. Tobenkin, and J. W. Roberts, “LQR-Trees: Feedback motion planning via sums-of-squares verification,” *The Int. J. of Rob. Res.*, vol. 29, no. 8, pp. 1038–1052, 2010.

[8] G. Goretin, A. Perez, R. Platt, and G. Konidaris, “Optimal sampling-based planning for linear-quadratic kinodynamic systems,” in *IEEE Internat. Conference on Robotics and Automation*, 2013, pp. 2429–2436.

[9] G. P. Kontoudis and K. G. Vamvoudakis, “Kinodynamic motion planning with continuous-time Q-learning: An online, model-free, and safe navigation framework,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 12, pp. 3803–3817, 2019.

[10] Y. Li, Z. Littlefield, and K. E. Bekris, “Asymptotically optimal sampling-based kinodynamic planning,” *The International Journal of Robotics Research*, vol. 35, no. 5, pp. 528–564, 2016.

[11] R. Shome and L. E. Kavraki, “Asymptotically optimal kinodynamic planning using bundles of edges,” in *2021 IEEE International Conference on Robotics and Automation*, 2021, pp. 9988–9994.

[12] J. D. Gammell and M. P. Strub, “Asymptotically optimal sampling-based motion planning methods,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 4, pp. 295–318, 2021.

[13] F. Atas, G. Cielniak, and L. Grimstad, “Benchmark of sampling-based optimizing planners for outdoor robot navigation,” in *International Conference on Intelligent Autonomous Systems*, 2023, pp. 231–243.

- [14] H.-T. L. Chiang, J. Hsu, M. Fiser, L. Tapia, and A. Faust, "RL-RRT: Kinodynamic motion planning via learning reachability estimators from RL policies," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 4298–4305, 2019.
- [15] C. Verginis, D. V. Dimarogonas, and L. Kavraki, "Sampling-based motion planning for uncertain high-dimensional systems via adaptive control," in *Algor. Foundations of Robotics XIV*, 2021, pp. 159–175.
- [16] P. Rousseas, C. Bechlioulis, and K. J. Kyriakopoulos, "Harmonic-based optimal motion planning in constrained workspaces using reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 2005–2011, 2021.
- [17] B. Kiumarsi, K. G. Vamvoudakis, H. Modares, and F. L. Lewis, "Optimal and autonomous control using reinforcement learning: A survey," *IEEE Tr. on Neur. Net. and Learn. Sys.*, vol. 29, no. 6, pp. 2042–2062, 2018.
- [18] P. Vlantis, C. Vrohidis, C. P. Bechlioulis, and K. J. Kyriakopoulos, "Robot navigation in complex workspaces using harmonic maps," in *IEEE Intern. Conf. on Robotics and Automation*, 2018, pp. 1726–1731.
- [19] D. Hsu, R. Kindel, J.-C. Latombe, and S. Rock, "Randomized kinodynamic motion planning with moving obstacles," *The International Journal of Robotics Research*, vol. 21, no. 3, pp. 233–255, 2002.
- [20] R. E. Allen and M. Pavone, "A real-time framework for kinodynamic planning in dynamic environments with application to quadrotor obstacle avoidance," *Rob. and Aut. Sys.*, vol. 115, pp. 174–193, 2019.
- [21] A. Majumdar and R. Tedrake, "Funnel libraries for real-time robust feedback motion planning," *The International Journal of Robotics Research*, vol. 36, no. 8, pp. 947–982, 2017.
- [22] D. Zheng and P. Tsiotras, "Sampling-based kinodynamic motion planning using a neural network controller," in *AIAA Scitech 2021 Forum*, 2021, p. 1754.
- [23] L. Li, Y. Miao, A. H. Qureshi, and M. C. Yip, "MPC-MPNet: Model-predictive motion planning networks for fast, near-optimal planning under kinodynamic constraints," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4496–4503, 2021.
- [24] M. K. M. Jaffar and M. Otte, "PiP-X: Funnel-based online feedback motion planning/replanning in dynamic environments," in *Workshop on the Algorithmic Foundations of Robotics*, 2022, pp. 132–148.
- [25] G. P. Kontoudis and K. G. Vamvoudakis, "Robust kinodynamic motion planning using model-free game-theoretic learning," in *American Control Conference*, 2019, pp. 273–278.
- [26] A. Kanellopoulos and K. G. Vamvoudakis, "Non-equilibrium dynamic games and cyber-physical security: A cognitive hierarchy approach," *Systems & Control Letters*, vol. 125, pp. 59–66, 2019.
- [27] K. G. Vamvoudakis and N.-M. T. Kokolakis, "Synchronous reinforcement learning-based control for cognitive autonomy," *Foundations and Trends® in Systems and Control*, vol. 8, no. 1–2, pp. 1–175, 2020.
- [28] N.-M. T. Kokolakis, A. Kanellopoulos, and K. G. Vamvoudakis, "Bounded rationality in differential games: A reinforcement learning-based approach," in *Handbook of Reinforcement Learning and Control*. Springer, 2021, pp. 467–489.
- [29] G. P. Kontoudis, Z. Xu, and K. G. Vamvoudakis, "Online, model-free motion planning in dynamic environments: An intermittent, finite horizon approach with continuous-time Q-learning," in *American Control Conference*, 2020, pp. 3873–3878.
- [30] J. Netter, G. P. Kontoudis, and K. G. Vamvoudakis, "Bounded rational RRT-QX: Multi-agent motion planning in dynamic human-like environments using cognitive hierarchy and Q-learning," in *IEEE Conference on Decision and Control*, 2021, pp. 3597–3602.
- [31] P. A. Ioannou and J. Sun, *Robust Adaptive Control*. Courier Corporation, 2012.
- [32] F. L. Lewis, D. Vrabie, and V. L. Syrmos, *Optimal Control*, 3rd ed. John Wiley & Sons, 2012.
- [33] D. Vrabie, O. Pastravanu, M. Abu-Khalaf, and F. L. Lewis, "Adaptive optimal control for continuous-time linear systems based on policy iteration," *Automatica*, vol. 45, no. 2, pp. 477–484, 2009.
- [34] K. G. Vamvoudakis, "Q-learning for continuous-time linear systems: A model-free infinite horizon optimal control approach," *Systems & Control Letters*, vol. 100, pp. 14–20, 2017.
- [35] Y. Yang, B. Kiumarsi, H. Modares, and C. Xu, "Model-free  $\lambda$ -policy iteration for discrete-time linear quadratic regulation," *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [36] P. Tabuada, "Event-triggered real-time scheduling of stabilizing control tasks," *IEEE Tr. on Aut. Con.*, vol. 52, no. 9, pp. 1680–1685, 2007.
- [37] K. G. Vamvoudakis and H. Ferraz, "Event-triggered H-infinity control for unknown continuous-time linear systems using Q-learning," in *IEEE Conference on Decision and Control*, 2016, pp. 1376–1381.
- [38] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [39] G. Chowdhary, T. Yucelen, M. Mühlegg, and E. N. Johnson, "Concurrent learning adaptive control of linear systems with exponentially convergent bounds," *Inte. Journal of Adaptive Control and Signal Processing*, vol. 27, no. 4, pp. 280–301, 2013.
- [40] K. G. Vamvoudakis, M. F. Miranda, and J. P. Hespanha, "Asymptotically stable adaptive-optimal control algorithm with saturating actuators and relaxed persistence of excitation," *IEEE Tr. on Neural Net. and Learning Sys.*, vol. 27, no. 11, pp. 2386–2398, 2016.
- [41] Y. Yang, Y. Pan, C.-Z. Xu, and D. C. Wunsch, "Hamiltonian-driven adaptive dynamic programming with efficient experience replay," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [42] G. P. Kontoudis, K. G. Vamvoudakis, and Z. Xu, "RRT-QX real-time kinodynamic motion planning in dynamic environments with continuous-time reinforcement learning," in *Brain and Cognitive Intelligence Control in Robotics*. CRC Press, 2022, pp. 1–19.
- [43] A. J. Van Der Schaft, "L<sub>2</sub>-gain analysis of nonlinear systems and nonlinear state-feedback H<sub>∞</sub> control," *IEEE Transactions on Automatic Control*, vol. 37, no. 6, pp. 770–784, 1992.
- [44] A. E. Bryson and Y.-C. Ho, *Applied Optimal Control: Optimization, Estimation, and Control (revised edition)*. Routledge, 1975.
- [45] T. Başar and P. Bernhard, *H<sub>∞</sub> Optimal Control and Related Minimax Design Problems: A Dynamic Game Approach*. Springer, 2008.
- [46] J. P. Hespanha, *Noncooperative Game Theory: An Introduction for Engineers and Computer Scientists*. Princeton University Press, 2017.
- [47] C. Berg, J. P. R. Christensen, and P. Ressel, *Harmonic Analysis on Semi-groups: Theory of Positive Definite and Related Functions*. Springer, 1984, vol. 100.
- [48] W. M. Haddad, V. Chellaboina, and S. G. Nersisov, *Impulsive and Hybrid Dynamical Systems: Stability, Dissipativity, and Control*. Princeton University Press, 2014, vol. 49.
- [49] J. P. Hespanha, D. Liberzon, and A. R. Teel, "Lyapunov conditions for input-to-state stability of impulsive systems," *Automatica*, vol. 44, no. 11, pp. 2735–2744, 2008.
- [50] M. De Berg, M. Van Kreveld, M. Overmars, and O. Schwarzkopf, *Computational Geometry*. Springer, 1997.
- [51] J. Zhang, K. H. Johansson, J. Lygeros, and S. Sastry, "Zeno hybrid systems," *International Journal of Robust and Nonlinear Control*, vol. 11, no. 5, pp. 435–451, 2001.
- [52] M. Otte, "RRT-X (dynamic obstacles)," [Online]. Available: [http://ottelab.com/html\\_stuff/code.html#RRTXcode](http://ottelab.com/html_stuff/code.html#RRTXcode), August 2023.
- [53] K. G. Vamvoudakis, "Event-triggered optimal adaptive control algorithm for continuous-time nonlinear systems," *IEEE/CAA Journal of Automatica Sinica*, vol. 1, no. 3, pp. 282–293, 2014.
- [54] H. Khalil, *Nonlinear Systems*. Prentice Hall, 2002.

## APPENDIX

*Proof of Theorem 2.* The proof follows from [53, Theorem 1]. We adopt the positive definite function  $\mathcal{V} := V^*$  as a Lyapunov candidate that satisfies the time-triggered HJI equation (9) with  $\mathcal{V}(0; \cdot) = 0$ . Then, the orbital derivative along the solution of (1) with the intermittent controller for  $t \in [r_j, r_{j+1})$  yields,

$$\dot{\mathcal{V}} = \frac{\partial V^*}{\partial \bar{x}} \dot{\bar{x}} + \frac{\partial V^*}{\partial t} = \frac{\partial V^*}{\partial \bar{x}} (A\bar{x} + Bu_d^* + Fd^*) + \frac{\partial V^*}{\partial t}. \quad (45)$$

After writing the time-triggered HJI equation (9) as,

$$-\frac{1}{2}(\bar{x}^\top M \bar{x} + u_c^{*\top} R u_c^* - \gamma^2 \|d^*\|^2) \equiv \frac{\partial V^*}{\partial \bar{x}} (A\bar{x} + Bu_c^* + Fd^*) + \frac{\partial V^*}{\partial t},$$

we rewrite the orbital derivative (45) as,

$$\begin{aligned} \dot{\mathcal{V}} &= -\frac{1}{2}\bar{x}^\top M \bar{x} - \frac{1}{2}u_c^{*\top} R u_c^* + \frac{1}{2}\gamma^2 \|d^*\|^2 + \frac{\partial V^*}{\partial \bar{x}} B(u_d^* - u_c^*) \\ &= -\frac{1}{2}\bar{x}^\top M \bar{x} + \frac{1}{2}(u_c^* - u_d^*)^\top R (u_c^* - u_d^*) - \frac{1}{2}u_d^{*\top} R u_d^* + \frac{1}{2}\gamma^2 \|d^*\|^2. \end{aligned}$$



Per Fact 1, we have  $(u_c^* - u_d^*)^\top R(u_c^* - u_d^*) \leq \bar{\lambda}(R)L(t)^2\|e\|^2$ . Thus, the orbital derivative (45) is upper-bounded by,

$$\dot{V} \leq -\frac{1}{2}(\beta^2\lambda(M) + (1 - \beta^2)\lambda(M))\|\bar{x}\|^2 + \frac{1}{2}\bar{\lambda}(R)L(t)^2\|e\|^2 - \frac{1}{2}\lambda(R)\|u_d^*\|^2 + \frac{1}{2}\gamma^2\|d^*\|^2. \quad (46)$$

Hence, if the following inequality is satisfied  $\forall t \in [r_j, r_{j+1})$ ,

$$\|e\|^2 \leq \frac{(1 - \beta^2)\lambda(M)\|\bar{x}\|^2 + \lambda(R)\|u_d^*\|^2 - \gamma^2\|d^*\|^2}{\bar{\lambda}(R)L(t)^2}, \quad (47)$$

then the right-hand side of (46) will be less than or equal to 0, so  $\dot{V} \leq 0$ , and the equilibrium point of the closed-loop system will be asymptotically stable as  $T \rightarrow \infty$ . ■

*Proof of Theorem 3.* We start by considering the flow dynamics (39). We take the Lyapunov candidate as,

$$\mathcal{L}(\psi; t) := V^*(\bar{x}; t) + V^*(\hat{x}; t) + \frac{1}{2}\|\tilde{W}_c\|^2 + \frac{1}{2}\text{tr}\{\tilde{W}_a^\top \tilde{W}_a\} + \frac{1}{2}\text{tr}\{\tilde{W}_d^\top \tilde{W}_d\} > 0, \quad \forall t \geq 0, \quad (48)$$

where  $\psi$  is the augmented state as defined in Section IV-D. The time derivative of (48) is partitioned to  $\dot{\mathcal{L}} = T_1 + T_2 + T_3 + T_4 + T_5$  with,

$$T_1 = \frac{1}{2}\bar{x}^\top \dot{P}(t)\bar{x} + \bar{x}^\top P(t)(A\bar{x} + B\hat{u}_d + F\hat{d}), \quad (49)$$

$$T_3 = \tilde{W}_c^\top \dot{\tilde{W}}_c = -\alpha_c \tilde{W}_c^\top \left( \frac{\sigma\sigma^\top}{(1 + \sigma^\top\sigma)^2} + \Lambda \right) \tilde{W}_c, \quad (50)$$

$$T_5 = \text{tr}\{\tilde{W}_d^\top \dot{\tilde{W}}_d\} = -\alpha_d \text{tr}\left\{ \tilde{W}_d^\top \frac{(\xi(t)\bar{x})(\xi(t)\bar{x})^\top}{1 + (\xi(t)\bar{x})^\top(\xi(t)\bar{x})} \tilde{W}_d + \tilde{W}_d^\top \frac{(\xi(t)\bar{x})\bar{x}^\top}{1 + (\xi(t)\bar{x})^\top(\xi(t)\bar{x})} \tilde{Q}_{dd} Q_{dd}^{-1} \right\}, \quad (51)$$

and  $T_2 = 0$ ,  $T_4 = 0$ . The terms  $T_2$  and  $T_4$  vanish because they are updated only at jumps and remain constant in the remaining time. We define  $\Xi := \frac{\sigma\sigma^\top}{(1 + \sigma^\top\sigma)^2} + \Lambda$  to upper bound  $T_3$  (50), i.e.,  $T_3 \leq -\alpha_c\lambda(\Xi)\|\tilde{W}_c\|^2$ . Since  $\frac{\sigma\sigma^\top}{(1 + \sigma^\top\sigma)^2} \succeq 0$  and  $\Lambda \succ 0$ , we know  $\lambda(\Xi) > 0$ . Thus, with  $\alpha_c \in \mathbb{R}^+$ , the following inequality holds,

$$T_3 \leq -\alpha_c\lambda(\Xi)\|\tilde{W}_c\|^2 \leq 0. \quad (52)$$

Next, the first term  $T_1$  (49) yields,

$$T_1 \leq -\frac{1}{2}\lambda(M)\|\bar{x}\|^2 + \frac{1}{2}\bar{\lambda}(R)\|\tilde{W}_a^\top \mu(t)\bar{x} - \hat{W}_a^\top \mu(t)e\|^2 - \frac{1}{2}\lambda(R)\|\hat{u}_d\|^2 - \frac{1}{2}\gamma^2\|\tilde{W}_d^\top \xi(t)\bar{x}\|^2 + \frac{1}{2}\gamma^2\|\hat{d}\|^2, \quad (53)$$

where  $\bar{x}^\top P(t)B = -u_c^{*\top}R$ , and  $\bar{x}^\top P(t)F = \gamma^2 d^{*\top}$ . Considering the actor weights remain constant during the flow dynamics, we have  $L_1(t) \mapsto \mathbb{R}^+$  a strictly positive-definite function [47] such that  $\|\tilde{W}_a^\top \mu(t)\| \leq L_1(t)$ . Thus, by adding and subtracting  $\frac{1}{2}\beta^2\lambda(M)\|\bar{x}\|^2$  to the right-hand side, and by using Fact 1, the Young's inequality, and (53), we obtain,

$$T_1 \leq -\frac{1}{2}((1 - \beta^2)\lambda(M) - 2\bar{\lambda}(R)L_1(t)^2)\|\bar{x}\|^2 - \frac{1}{2}\beta^2\lambda(M)\|\bar{x}\|^2 + 2\bar{\lambda}(R)(L(t)^2 + L_1(t)^2)\|e\|^2 - \frac{1}{2}\lambda(R)\|\hat{u}_d\|^2 - \frac{1}{2}\gamma^2\|\tilde{W}_d^\top \xi(t)\bar{x}\|^2 + \frac{1}{2}\gamma^2\|\hat{d}\|^2. \quad (54)$$

Provided that in the time interval  $t \in [r_j, r_{j+1})$ , the upper

bound inequality (47) holds, then (54) yields,

$$T_1 \leq -\frac{1}{2}(\beta^2\lambda(M) - 2\bar{\lambda}(R)L_1(t)^2)\|\bar{x}\|^2 - \frac{1}{2}\gamma^2\|\tilde{W}_d^\top \xi(t)\bar{x}\|^2. \quad (55)$$

Now consider the last term  $T_5$  (51). It yields,

$$T_5 \leq -\frac{\alpha_d}{2} \frac{\|\bar{x}^\top \xi(t)^\top \tilde{W}_d\|^2}{1 + (\xi(t)\bar{x})^\top(\xi(t)\bar{x})} + \frac{\alpha_d}{2} \eta \bar{\lambda} \left( \frac{\xi(t)Q_{dd}^{-1}}{1 + (\xi(t)\bar{x})^\top(\xi(t)\bar{x})} \right) \|\bar{x}\|^2, \quad (56)$$

where  $\eta$  is a constant of unity order. Given the upper bounds (52), (55) and (56), we have,

$$\dot{\mathcal{L}}(\psi; t) \leq -\frac{1}{2}(\beta^2\lambda(M) - 2\bar{\lambda}(R)L_1(t)^2 - \alpha_d \eta \bar{\lambda} \left( \frac{\xi(t)Q_{dd}^{-1}}{1 + (\xi(t)\bar{x})^\top(\xi(t)\bar{x})} \right) \|\bar{x}\|^2 - \frac{1}{2}\gamma^2\|\tilde{W}_d^\top \xi(t)\bar{x}\|^2 - \frac{\alpha_d}{2} \frac{\|\tilde{W}_d^\top \xi(t)\bar{x}\|^2}{1 + (\xi(t)\bar{x})^\top(\xi(t)\bar{x})} - \alpha_c\lambda(\Xi)\|\tilde{W}_c\|^2).$$

Thus, if the following inequalities are satisfied,

$$0 < \alpha_d \leq \frac{\beta^2\lambda(M) - 2\bar{\lambda}(R)L_1(t)^2}{\eta \bar{\lambda} \left( \frac{\xi(t)Q_{dd}^{-1}}{1 + (\xi(t)\bar{x})^\top(\xi(t)\bar{x})} \right)}, \quad \frac{\lambda(M)}{\bar{\lambda}(R)} > \frac{2L_1(t)^2}{\beta^2},$$

then  $\dot{\mathcal{L}}(\psi; t)$  is non-positive for all  $\psi$  and  $t \geq t_0$ . By defining  $W_1(\psi; t) = W_2(\psi; t) := V^*(\bar{x}; t) + \frac{1}{2}\|\tilde{W}_c\|^2 + \frac{1}{2}\text{tr}\{\tilde{W}_d^\top \tilde{W}_d\} > 0$ ,  $W_1(\psi; t) \leq \mathcal{L}(\psi; t) \leq W_2(\psi; t)$  is satisfied. Hence, the origin  $\psi_e = 0$  is uniformly stable per the Lyapunov stability theorem. Also, given that  $\mathcal{L}(\psi; t)$  is lower-bounded, non-increasing, and that its time derivative  $\dot{\mathcal{L}}(\psi; t) = T_1 + T_3 + T_5$  is also bounded, the Lyapunov function in (48) is uniformly continuous, which satisfies the Barbalat's lemma,  $\mathcal{L}(\psi; t) \rightarrow 0$  as  $t \rightarrow \infty$ . Since  $\dot{\mathcal{L}}(\psi; t) < 0$ , asymptotic stability holds from the Lyapunov stability theorem. Moreover,  $W_1(\psi; t)$  is radially unbounded with respect to  $\|\bar{x}\|$ ,  $\|\tilde{W}_c\|$ ,  $\|\tilde{W}_d^\top \xi(t)\bar{x}\|^2$ , and thus globally properties hold. Therefore, the equilibrium point of the flow dynamics at the origin  $\psi_e = 0$  is globally uniformly asymptotically stable as  $T \rightarrow \infty$  [54].

We continue with the jump dynamics (40) comprising the sampled states and the control actor weight updates. Let us define the Lyapunov function at  $t = r_j$  as,

$$\Delta\mathcal{L}(\psi; t) := \Delta V_1(\bar{x}^+, \bar{x}(r_j); t) + \underbrace{(V^*(\hat{x}^+; t) - V^*(\hat{x}(r_j); t))}_{\Delta V_2} + \Delta V_3(\tilde{W}_c^+, \tilde{W}_c(r_j)) + \Delta V_5(\tilde{W}_d^+, \tilde{W}_d(r_j)) + \underbrace{\frac{1}{2\alpha_a}(\text{tr}\{\tilde{W}_a^{+\top} \tilde{W}_a^+\} - \text{tr}\{\tilde{W}_a(r_j)^\top \tilde{W}_a(r_j)\})}_{\Delta V_4}, \quad (57)$$

where  $\Delta\mathcal{L}(\psi; t) > 0$ . Note that  $\bar{x}$ ,  $\tilde{W}_c$ ,  $\tilde{W}_d$  are continuous even at the triggering events, hence they all vanish during the jump dynamics, i.e.,  $\Delta V_1(\bar{x}^+, \bar{x}(r_j); t) = 0$ ,  $\Delta V_3(\tilde{W}_c^+, \tilde{W}_c(r_j)) = 0$ , and  $\Delta V_5(\tilde{W}_d^+, \tilde{W}_d(r_j)) = 0$ . With  $\hat{x}^+ = \bar{x}(r_j)$  and  $\hat{x}(r_j) = \bar{x}(r_{j-1})$  at the jump, the second term of (57) yields,

$$\Delta V_2 = V^*(\hat{x}^+) - V^*(\hat{x}(r_j)) = V^*(\bar{x}(r_j)) - V^*(\bar{x}(r_{j-1})) \leq 0,$$

because of the convergence of  $\bar{x}$ . Hence, as  $T \rightarrow \infty$ , the sampled state asymptotically converges to the origin, i.e.,

$\|\hat{x}(r_j)\| \rightarrow 0$ .

Using  $Q_{u_d u_d} = R$  and Young's inequality, taking entry-wise norms for the fourth term of (57) results in,

$$\begin{aligned} \Delta V_4 \leq & -\|\tilde{W}_a^\top \mu(t) \bar{x}(r_j)\|^2 \left(1 - \frac{1 + \alpha_a}{4\bar{\lambda}(R)} - \frac{\alpha_a}{8}\right) \\ & + \frac{1 + \alpha_a}{4\bar{\lambda}(R)} \|\tilde{Q}_{xu_d}\|^2 + \frac{\alpha_a}{2\bar{\lambda}(R)^2} \|\tilde{Q}_{xu_d}\|^2. \end{aligned} \quad (58)$$

Therefore, the inequality  $\Delta V_4 < 0$  is satisfied if we guarantee in (58) that  $\tilde{W}_a$  remains outside the compact

$$\text{set } \Omega = \left\{ \tilde{W}_a \in \mathbb{R}^{n \times m} \mid \|\tilde{W}_a\| \leq \sqrt{\frac{\frac{1 + \alpha_a}{4\bar{\lambda}(R)} \|\tilde{Q}_{xu_d}\|^2 + \frac{\alpha_a}{2\bar{\lambda}(R)^2} \|\tilde{Q}_{xu_d}\|^2}{1 - \frac{1 + \alpha_a}{4\bar{\lambda}(R)} - \frac{\alpha_a}{8}}} \right\}.$$

The last two constraints arise from the denominator,

$$0 < \alpha_a < \frac{8\bar{\lambda}(R) - 2}{\bar{\lambda}(R) + 2}, \quad \bar{\lambda}(R) > 0.25.$$

Since the elements in  $\Omega$  are asymptotically stable for  $T \rightarrow \infty$ , the set becomes a single point, thus  $\|\tilde{W}_a\| \rightarrow 0$ .  $\blacksquare$