

# Assured Learning-Based Optimal Control subject to Timed Temporal Logic Constraints

Filippos Fotiadis<sup>1</sup>, Christos K. Verginis<sup>2</sup>, Kyriakos G. Vamvoudakis<sup>1</sup>, and Ufuk Topcu<sup>2</sup>

**Abstract**—We develop an algorithm for the optimal control of systems governed by unknown, nonlinear dynamics to deliver tasks expressed as timed temporal logic constraints. The algorithm first computes a sequence of points—along with associated time stamps—in the operating environment such that, if the system follows the sequence, it completes its task. For the second step of the algorithm, we develop a novel data-driven, on-the-fly control mechanism that learns how to transition from a point in the sequence to the next within a pre-specified time horizon accounting for the unknown dynamics, unsafe zones in the operating environment and additional optimality criteria. We show that, after a finite period of data gathering, the resulting controller guarantees that the system indeed follows the sequence of points, leading to the satisfaction of the task.

## I. INTRODUCTION

Autonomous systems are prone to failures and abrupt changes that might render the underlying dynamics unknown. Hence, the control of such systems necessitates data-driven, learning-based techniques. Moreover, abrupt changes in the dynamics prevent the use of data obtained offline, and the learning-based techniques should rely on data obtained on the fly from the current system trajectory. Additionally, expressing the objectives of autonomous systems via temporal logic languages has gained significant attention recently, since temporal logic can describe more complex tasks than the well-studied point-to-point navigation [1]. A special form of temporal logic, namely *timed temporal logic*, offers the incorporation of time constraints in the planning objectives, providing a rich variety of tasks [2]. At the same time, resource limitations call for algorithms that minimize the exerted control effort of the underlying system by solving the optimal control problem [3].

This paper addresses the optimal control problem of an *unknown* control-affine system to deliver tasks expressed as timed temporal logic constraints. The system is assumed to be continuous in state and time and operating in an environment with unsafe zones. Our contribution with respect to the related literature lies in the integration of timed temporal tasks with control optimality for unknown nonlinear systems.

We develop a two-step algorithm to solve the aforementioned problem. The first step is the computation of a discrete

timed path, i.e., a sequence of points to be visited at specific time stamps, that yields the execution of the task if followed by the system. The second step, which constitutes the main contribution of our work, is the design of a control algorithm that exhibits the following properties: (i) it achieves the sequential navigation of the system to the points dictated by the computed path in the given time stamps, (ii) it minimizes the exerted control effort, and (iii) it guarantees the avoidance of the unsafe zones. In particular, we transform the problem to a finite-horizon optimal control problem with safety constraints and we use data obtained online from the current trajectory to accommodate the unknown dynamics. We prove that, after obtaining a sufficient amount of data, the system learns to navigate among the predefined points within the time intervals dictated by the derived path while minimizing the control effort and avoiding the unsafe zones, which leads to the successful execution of the task.

There exist numerous related works that consider planning and control under timed temporal logic specifications [4]–[15]. Most of the aforementioned works, however, consider simplistic single integrator models [5], [7], [11], finite-state systems [15] or neglect entirely the underlying dynamics [4]. The works [8]–[10], [12]–[14] consider more complex models that are either fully [10], [12]–[14] or partially [8], [9] known; [2] assumes unknown dynamics, restricted, however, to second-order Lagrangian models with positive-definite input matrices. This paper considers a larger class of systems, governed by fully unknown nonlinear dynamics.

Another issue with the related works on timed temporal logic-based planning is the lack of optimality characteristics; [11] aims to minimize the system’s control effort by online modifying the timed paths, whereas [8] embeds a policy improvement algorithm to a feedback control law for simultaneous satisfaction of timed temporal specifications and minimization of a given cost. The work in [13] considers the optimal control problem by incorporating timed temporal logic specifications as constraints using quadratic programs. Nevertheless, except for using the underlying dynamics, the aforementioned works fail to guarantee optimality of the resulting controller. In this paper, on the other hand, we propose a neural-network-based learning scheme that guarantees the optimality of the resulting controller up to an approximation error that depends on the size of the neural network. In particular, we extend previous works on actor-critic learning [16]–[18] by solving a series of optimal control problems over finite time horizons for the safe timed transitions among the predefined points that are related to the timed temporal task. We provide formal guarantees

<sup>1</sup>F. Fotiadis and K. G. Vamvoudakis are with the Georgia Institute of Technology, Atlanta, GA, USA. Email: {ffotiadis, kyriakos}@gatech.edu. <sup>2</sup>C. K. Verginis and U. Topcu are with The University of Texas at Austin, Austin, TX, USA. Email: {cverginis, utopcu}@utexas.edu. This work was supported in part, by ARO under grant No. W911NF-19-1-0270, by ONR Minerva under grant No. N00014-18-1-2160, and by NSF under grant Nos. CAREER CPS-1851588 and SATC-1801611, by NASA ULI under grant number 80NSSC20M0161 and by the Onassis Foundation-Scholarship ID: F ZQ 064 – 1/2020 – 2021.

regarding the optimality of the resulting closed-loop system, which, according to the authors' best knowledge, has not been considered before for the finite-horizon optimal control problem with unknown continuous-time nonlinear dynamics.

## II. PRELIMINARIES

*Notation:* We denote by  $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$  the set of nonnegative integer numbers, where  $\mathbb{N}$  is the set of natural numbers. The sets of  $n$ -dimensional nonnegative and positive reals, with  $n \in \mathbb{N}$ , are denoted by  $\mathbb{R}_{\geq 0}^n$  and  $\mathbb{R}_{> 0}^n$ , respectively;  $Z_1 \otimes Z_2$  is the Kronecker product of matrices  $Z_1$  and  $Z_2$ . The operator  $\text{vec}(\cdot)$  denotes the vectorization of a matrix. Given an infinite sequence  $s = s_0 s_1 s_2 \dots$ , we denote its  $j$ -suffix by  $\text{suff}(s, j) = s_j s_{j+1} \dots$ , respectively;  $I_q \in \mathbb{R}^{q \times q}$  denotes the identity matrix. The closed ball centered at  $c_k$  with radius  $r_k$  is denoted by  $\bar{B}(c_k, r_k)$ .

**Definition 1** ([19]). A time sequence  $t_1, t_2, \dots$  is a (infinite unless otherwise stated) sequence of time values  $t_j \in \mathbb{R}_{\geq 0}$ , for all  $j \in \mathbb{N}$ , satisfying (i)  $t_j < t_{j+1}$ , for all  $j \in \mathbb{N}$  and (ii) for all  $t' \in \mathbb{R}_{\geq 0}$  there exists  $j \geq 1$  such that  $t_j \geq t'$ .  $\square$

An *atomic proposition* is a statement over the variables or parameters of a problem that is either True ( $\top$ ) or False ( $\perp$ ), and let  $\mathcal{AP}$  be a finite set of such atomic propositions.

**Definition 2.** Let  $\mathcal{AP}$  be a finite set of atomic propositions. A *timed word*  $w$  over  $\mathcal{AP}$  is an infinite sequence  $w := (w_1, t_1)(w_2, t_2) \dots$  where  $w_1, w_2, \dots$  is an infinite word over  $2^{\mathcal{AP}}$  and  $t_1, t_2, \dots$  is a time sequence.  $\square$

**Definition 3.** A *Weighted Transition System (WTS)* is a tuple  $(\Pi, S\Pi_0, \longrightarrow, \mathcal{AP}, \mathcal{L}, \gamma)$ , where  $\Pi$  is a finite set of states,  $\Pi_0 \subseteq \Pi$  is a set of initial states,  $\longrightarrow \subseteq S \times S$  is a transition relation,  $\mathcal{AP}$  is a finite set of atomic propositions,  $\mathcal{L} : \Pi \rightarrow 2^{\mathcal{AP}}$  is a labeling function, and  $\gamma : \longrightarrow \rightarrow \mathbb{R}_{> 0}$  is a map that assigns a positive weight to each transition.  $\square$

**Definition 4.** A *timed run* of a WTS is an infinite sequence  $r = (r_1, t_1)(r_2, t_2) \dots$ , such that  $r_1 \in S_0$  and  $r_j \in S$ ,  $(r_j, r_{j+1}) \in \longrightarrow$ , for all  $j \in \mathbb{N}$ . The time stamps  $t_j$  are inductively defined with  $t_1 = 0$  and  $t_{j+1} = t_j + \gamma(r_j, r_{j+1})$ , for all  $j \in \mathbb{N}$ . The timed run  $r$  generates a timed word  $w(r) = w_1(r_1), w_2(r_2), \dots = (\mathcal{L}(r_1), t_1), (\mathcal{L}(r_2), t_2), \dots$  over the set  $2^{\mathcal{AP}}$ , where  $\mathcal{L}(r_j)$  is the subset of atomic propositions  $\mathcal{AP}$  that are true at state  $r_j$  at time  $t_j$ ,  $j \in \mathbb{N}$ .  $\square$

The syntax of a timed temporal logic formula over  $\mathcal{AP}$  is defined by a grammar that has the form

$$\varphi := \text{p} \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \bigcirc_I \varphi \mid \diamond_I \varphi \mid \square_I \varphi \mid \varphi_1 \mathcal{U}_I \varphi_2, \quad (1)$$

where  $\varphi \in \mathcal{AP}$ , and  $\bigcirc$ ,  $\diamond$ ,  $\square$ , and  $\mathcal{U}$  are the next, future, always, and until operators, respectively;  $I$  is a nonempty time interval in one of the followings forms:  $[i_1, i_2], [i_1, i_2), (i_1, i_2], (i_1, i_2), [i_1, \infty), (i_1, \infty)$ , with  $i_1, i_2 \in \mathbb{Q}$ . Several languages are subsets of the form (1), such as Metric Temporal Logic (MTL), Metric Interval Temporal Logic (MITL), Bounded MTL, coFlat MTL, or Time Window Temporal Logic (TWTL) [20], [21]. Here we define

the generalized semantics of (1) over discrete observations (point-wise semantics) [22]. The next definition considers the satisfaction of a formula by a timed run.

**Definition 5.** [22], [23] Given a sequence  $R = (\pi_0, t_0)(\pi_1, t_1) \dots$  and a timed formula  $\varphi$ , we define  $(R, i) \models \varphi$ ,  $i \in \mathbb{N}_0$  ( $R$  satisfies  $\varphi$  at  $i$ ) as follows:

$$\begin{aligned} (R, i) \models p &\Leftrightarrow p \in \mathcal{L}(\pi_i), & (R, i) \models \neg\varphi &\Leftrightarrow (R, i) \not\models \varphi, \\ (R, i) \models \varphi_1 \wedge \varphi_2 &\Leftrightarrow (R, i) \models \varphi_1 \text{ and } (R, j) \models \varphi_2, \\ (R, i) \models \bigcirc_I \varphi &\Leftrightarrow (R, i+1) \models \varphi \text{ and } t_{j+i} - t_i \in I, \\ (R, i) \models \varphi_1 \mathcal{U}_I \varphi_2 &\Leftrightarrow \exists k \geq i \text{ such that } (R, k) \models \varphi_2, \\ & t_k - t_i \in I \text{ and } (R, m) \models \varphi_1, \forall m \in \{j, \dots, k\}. \end{aligned}$$

Also,  $\diamond_I \varphi = \top \mathcal{U}_I \varphi$  and  $\square_I \varphi = \neg \diamond_I \neg \varphi$ . Finally,  $R$  satisfies  $\varphi$ , denoted by  $R \models \varphi$ , if and only if  $(R, 0) \models \varphi$ .  $\square$

## III. PROBLEM FORMULATION

Consider, for all  $t \geq t_0 \geq 0$ , a nonlinear system with dynamics

$$\dot{x}(t) = f(x(t)) + g(x(t))u(x(t), t), \quad x(t_0) = x_0, \quad (2)$$

where  $x : [t_0, \infty) \rightarrow \mathbb{R}^n$  denotes the system's states with initial condition  $x_0 \in \mathbb{R}^n$  at  $t = t_0$ ,  $u : \mathbb{R}^n \times [t_0, \infty) \rightarrow \mathbb{R}^m$  is a control input, and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $g : \mathbb{R}^m \rightarrow \mathbb{R}^{n \times m}$  are unknown, locally Lipschitz functions.

Moreover, consider  $K \in \mathbb{N}$  points of interest in the state space, denoted by  $c_k \in \mathbb{R}^n$ , for  $k \in \mathcal{K} := \{1, \dots, K\}$ , and let  $\Pi := \{c_1, \dots, c_K\}$ . Each point  $c_k$ ,  $k \in \mathcal{K}$ , corresponds to certain properties of interest, which are expressed as Boolean variables via the finite set of atomic propositions  $\mathcal{AP}$ . The properties satisfied at each point are provided by the labeling function  $\mathcal{L} : \Pi \rightarrow 2^{\mathcal{AP}}$ . Informally,  $\mathcal{L}$  assigns to each point  $c_k$ ,  $k \in \mathcal{K}$ , the subset of the atomic propositions that hold true in that point. Since the aforementioned properties shared by a point of interest are naturally inherited to some neighborhood of that point, we also define for each  $k \in \mathcal{K}$  the region of interest  $\pi_k$ , corresponding to the point of interest  $c_k$ , as the set  $\pi_k := \bar{B}(c_k, \rho_k)$ , with  $\rho_k > 0$  chosen such that  $\pi_k \cap \pi_{k'} = \emptyset$ , for all  $k, k' \in \mathcal{K}$  with  $k \neq k'$ . The system is assumed to be in a  $\pi_k$  simply when  $x \in \pi_k$ . We further need the following assumption.

**Assumption 1.** It holds that  $f(c_k) = 0$  for all  $k \in \mathcal{K}$ .  $\square$

Along with  $\Pi$ , we further consider a set of  $K_o$  unsafe pairwise disjoint spherical zones  $\mathcal{O} := \{o_1, \dots, o_{K_o}\}$ , with  $o_k := \bar{B}(c_{o_k}, \rho_{o_k})$ ,  $k \in \mathcal{K}_o := \{1, \dots, K_o\}$  satisfying  $o_k \cap \pi_{k'} = \emptyset$ , for all  $(k, k') \in \mathcal{K}_o \times \mathcal{K}$ , which defines the free space  $\mathcal{F} := \mathbb{R}^n \setminus \mathcal{O}$ . We are interested in achieving timed temporal specifications over the atomic propositions  $\mathcal{AP}$  while avoiding the unsafe zones. We achieve that by guaranteeing safe timed transitions between the regions of interest in  $\Pi$ . We first need the following definition regarding the *behavior* of the system.

**Definition 6.** Consider an agent trajectory  $x : [t_0, \infty) \rightarrow \mathbb{R}^n$  of (2). Then, a *timed behavior* of  $x$  is the infinite sequence  $\text{b}(t_0) := (x(t_0), \sigma_0, t_0)(x(t_1), \sigma_1, t_1) \dots$ , where  $t_0, t_1, \dots$  is

a time sequence according to Definition 1,  $x(t_i) \in \pi_{j_i}$ ,  $j_i \in \mathcal{K}$  for all  $i \in \mathbb{N}_0$ , and  $\sigma_i = \mathcal{L}(\pi_{j_i}) \subseteq 2^{\mathcal{AP}}$ , i.e., the subset of atomic propositions that are true when  $x(t_j) \in \pi_{j_i}$ , for  $i \in \mathbb{N}_0$ . The timed behavior  $\mathbf{b}$  satisfies a timed formula  $\varphi$  *safely* if  $\mathbf{b}_\sigma(t_0) := (\sigma_0, t_0)(\sigma_1, t_1) \dots \models \varphi$  and  $x(t) \in \mathcal{F}$ , for all  $t \geq t_0$ . It *eventually* satisfies  $\varphi$  *safely* if there exists  $j \in \mathbb{N}$  such that  $\text{suff}(\mathbf{b}_\sigma(t_0), j) = \text{suff}(a_\sigma(t_0), j)$ , for some  $a_\sigma(t_0) \models \varphi$  and  $x(t) \in \mathcal{F}$ , for all  $t \geq t_j$ .  $\square$

We develop a learning-based control strategy such that the system learns how to safely execute transitions in  $\Pi$ , resulting in eventual satisfaction of  $\varphi$ , while also achieving optimality with respect to some user-defined cost. Note that eventual satisfaction implies that  $\varphi$  dictates repetitive tasks and/or tasks over long time horizons that the system is able to learn to execute. The latter, however, is not a restrictive assumption, since such tasks encompass the full potential of timed temporal logic languages.

Define now, for each point of interest  $c_i$ , the error  $e_i := x - c_i$ , evolving according to the dynamics

$$\dot{e}_i = F_i(e_i) + G_i(e_i)u := f(e_i + c_i) + g(e_i + c_i)u, \quad (3)$$

for all  $i \in \mathcal{K}$ , and define the performance criteria:

$$J(e_i(t_0), t_0, t_f, u) := \int_{t_0}^{t_f} r(e_i(\tau), u(e_i(\tau), \tau))d\tau, \quad (4)$$

with  $t_0 \geq 0$ ,  $t_f > t_0$ , where  $r(e, u) := Q(e) + S(u)$  is a metric of performance, with  $S(u) := u^T R u$ ,  $R > 0$ , and  $Q : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$  being a positive-definite function. This gives rise to the *timed behavior cost* of a timed behavior  $\mathbf{b}$ .

**Definition 7.** Consider a system closed-loop trajectory  $x : [t_0, \infty) \rightarrow \mathbb{R}^n$  along the control input  $u$  and the associated timed behavior  $\mathbf{b} = (x(t_0), \sigma_0, t_0)(x(t_1), \sigma_1, t_1) \dots$ , with  $x(t_i) \in \pi_{j_i}$ ,  $j_i \in \mathcal{K}$  for all  $i \in \mathbb{N}_0$ . The *timed behavior cost*  $J$  is the infinite sequence of functions  $J := J_0 J_1 \dots$ , where  $J_i := J(e_{j_{i+1}}(t_i), t_i, t_{i+1}, u)$ , for all  $i \in \mathbb{N}_0$ .  $\square$

The cost of the timed behavior naturally leads to the  $\epsilon$ -*optimal timed behavior* defined next, where  $\mathcal{A}(t_a, t_b)$  is the set of all functions from  $\mathbb{R}^n \times [t_a, t_b]$  to  $\mathbb{R}^m$ ,  $t_b > t_a \geq t_0$ :

**Definition 8.** Consider a system trajectory  $x : [t_0, \infty) \rightarrow \mathbb{R}^n$ . Given  $\epsilon > 0$ , its timed behavior  $\mathbf{b}(t_0) = (x(t_0), \sigma_0, t_0)(x(t_1), \sigma_1, t_1) \dots$  is said to be  $\epsilon$ -*optimal*, if the associated timed behavior cost  $J = J_0 J_1 \dots$  satisfies  $\|J_i - J_i^*\| \leq \epsilon$ , for all  $i \in \mathbb{N}$ , where  $J_i^* := \min_{\alpha \in \mathcal{A}(t_i, t_{i+1})} J(e_{j_{i+1}}(t_i), t_i, t_{i+1}, \alpha)$ .  $\square$

We can now state the problem considered in this work.

**Problem 1.** Let a system evolve with unknown dynamics (2), with initial position  $x(t_0) \in \pi_{j_0}$ ,  $j_0 \in \mathcal{K}$ . Given a timed formula  $\varphi$  over  $\mathcal{AP}$  and a labeling function  $\mathcal{L}$ , design a control law  $u : \mathbb{R}^n \times [t_0, \infty) \rightarrow \mathbb{R}^m$  that results in a solution  $x : [t_0, \infty) \rightarrow \mathbb{R}^n$ , which achieves an  $\epsilon$ -*optimal* timed behavior that eventually satisfies  $\varphi$  *safely*.  $\square$

The next sections describe our two-layered solution to Problem 1. We first synthesize a high-level timed path over

$\Pi$  that satisfies  $\varphi$ , by neglecting the unknown dynamics (2). Then, we design a novel learning-based control algorithm that learns how to execute safe timed transitions over  $\Pi$  based on data obtained online from the current trajectory, which leads to the eventual safe satisfaction of  $\varphi$ .

#### IV. HIGH-LEVEL PLAN GENERATION

The first ingredient of the proposed solution is the derivation of a high-level plan that satisfies the given formula  $\varphi$ . To this end, we abstract the motion of the system as a finite weighted transition system [1]

$$\mathcal{T} := (\Pi, \Pi_0, \longrightarrow, \mathcal{AP}, \mathcal{L}, \gamma) \quad (5)$$

where  $\Pi$  is the discretized state space,  $\Pi_0 \subseteq \Pi$  is the initial region, computed as  $\Pi_0 := \pi_{k_0}$ ,  $k_0 := \arg \min_{k \in \mathcal{K}} \{\|x(t_0) - c_k\|\}$ ,  $\longrightarrow \subseteq \Pi \times \Pi$  is a transition relation,  $\mathcal{AP}$  and  $\mathcal{L}$  are the set of atomic propositions and labeling function, respectively, defined in the previous section, and  $\gamma : \longrightarrow \rightarrow \mathbb{R}_{>0}$  is a map that assigns a positive weight to each transition. For now we assume that the system can execute the transitions among the regions in  $\Pi$  within the time interval dictated by  $\gamma$ ; the latter can be chosen according to several criteria, such as input capabilities of the system, Euclidean distance among points of interest, etc<sup>1</sup>. In the next section we will consider the control design for the execution of these timed transitions.

Given the transition system  $\mathcal{T}$  and the formula  $\varphi$ , we can apply standard formal verification methodologies in order to compute a timed path over  $\Pi$  that satisfies  $\varphi$ . The most common practice to achieve this is the following: Firstly,  $\varphi$  is algorithmically translated to a Timed Büchi Automaton (TBA)  $\mathcal{A}_B$ , a system consisting of a discrete set of states associated with  $\mathcal{AP}$ , whose accepting runs satisfy  $\varphi$  [1]. Secondly, we compute the product between the two discrete systems  $\tilde{\mathcal{T}} := \mathcal{T} \otimes \mathcal{A}_B$ ; Finally,  $\tilde{\mathcal{T}}$  is viewed as a graph and standard graph-based algorithms are used to derive a *timed* path that satisfies  $\varphi$ . This path has the prefix-suffix form

$$\mathbf{p} = (\pi_{k_0}, t_0) \dots (\pi_{k_{\mu-1}}, t_{k_{\mu-1}}) \left[ (\pi_{k_\mu}, t_{k_\mu}) \dots (\pi_{k_{\mu+\nu}}, t_{k_{\mu+\nu}}) \right]^\omega,$$

where  $\mu_1 := \mu + \nu$ ,  $\mu_2 := \mu + (\nu + 1)$ , for positive  $\mu$ ,  $\nu$ , where the superscript  $\omega$  denotes infinite repetition and  $\iota = 0, 1, \dots$  denotes the repetition index. The execution of  $\mathbf{p}$  produces a trajectory  $x(t)$ ,  $t \geq t_0$ , with timed behavior  $\mathbf{b}(t_0)$  that satisfies  $\varphi$ , i.e.,  $\mathbf{b}_\sigma(t_0) \models \varphi$  (see Definition 6). One can also obtain a timed path  $\mathbf{p}$  satisfying  $\varphi$  using optimization methodologies. In particular, it has been shown that the satisfaction of a timed temporal formula can be formulated as a Mixed Integer Linear Programming (MILP) problem [4], where binary variables are introduced to represent the several atomic propositions and the time constraints involved in  $\varphi$ .

After obtaining the timed path  $\mathbf{p}$ , we design in the next section a data-based learning control protocol that learns over time how to successfully execute the timed transitions in  $\mathbf{p}$  while avoiding the unsafe zones, leading to the eventual satisfaction of  $\varphi$ , as per Def. 6.

<sup>1</sup>One can also consider online reconfiguration algorithms that give an optimal time duration based on exerted control effort [11].

## V. OPTIMAL TRANSITION

This section describes the data-based optimal control design for the *optimal* timed transition among two regions  $\pi_k$ , and  $\pi_\ell$ , which is defined as follows.

**Definition 9.** Assume that  $x(t_k) \in \pi_k$ , for a  $t_k \in \mathbb{R}_{\geq 0}$ . Then, the system performs an *optimal* timed transition to  $\pi_\ell$ ,  $\ell \in \mathcal{K} \setminus \{k\}$ , denoted by  $\pi_k \longrightarrow$ , if it applies a time-varying feedback control law  $u : \mathbb{R}^n \times [t_k, t_\ell] \rightarrow \mathbb{R}^m$  such that, for some  $\delta \in \mathbb{R}_{> 0}$ , the solution of the closed loop system (2) satisfies the following:

- 1)  $x(t) \in \pi_\ell$  for all  $t \in [t_k + \delta, t_\ell]$ ,
- 2)  $x(t) \in \mathcal{F} \setminus \bigcup_{m \in \mathcal{K} \setminus \{\ell\}} \pi_m$  for all  $t \in [t_k, t_\ell]$ ,
- 3)  $u(x(t), t) = \arg J_i^*$ , where  $J_i^* = \min_{\alpha \in \mathcal{A}(t_k, t_\ell)} J(e_\ell(t_k), t_k, t_\ell, \alpha)$ .

The timed transition is  $\epsilon$ -optimal, if 3) is replaced by  $\|J(e_\ell(t_k), t_k, t_\ell, u(x(t), t)) - J_i^*\| \leq \epsilon$ .  $\square$

### A. Optimal Control with Soft Constrains

Evidently, it is not necessary that a control law  $u$  that minimizes (4) can always achieve the timed behavior described in Problem 1. Hence, the minimization of (4) has to be subject to some hard constraints imposing the desired timed behavior, or the desired behavior can be incorporated as a soft constraint in the cost (4). To achieve the latter, let us consider two regions of interest  $\pi_k$ ,  $\pi_\ell$  corresponding to two subsequent time instances  $t_k$  and  $t_\ell$ , with  $k, \ell \in \mathcal{K}$ . Then, we redefine the performance criterion (4) into:

$$J_\ell(e_\ell(t_k), t_k, u) := \int_{t_k}^{t_\ell} \left( \gamma r(e_\ell(\tau), u(e_\ell(\tau))) + L_{k,\ell}(e_\ell(\tau)) \right) d\tau + \phi(e_\ell(t_\ell)), \quad (6)$$

$$\text{subject to: } \dot{e}_\ell = F_\ell(e_\ell) + G_\ell(e_\ell)u, \quad (7)$$

where we also drop the final time  $t_\ell$  argument for brevity. In (6),  $\phi : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$  is a positive-definite term that penalizes the deviation of the terminal state  $e_\ell(t_\ell)$  from the point of interest  $c_\ell$ . In addition,  $L_{k,\ell} : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$  is another penalty term satisfying  $L_{k,\ell}(0) = 0$  and designed so that the system, under the controller that minimizes (6), avoids all unsafe zones  $\mathcal{O}$  and regions  $\pi_i$ ,  $i \in \mathcal{K} \setminus \{\ell\}$ . A possible implementation of  $L_{k,\ell}$  is given in the Appendix. Finally,  $\gamma > 0$  dictates a trade-off between a) ensuring avoidance of the unsafe zones and the regions  $\pi_i$ ,  $i \in \mathcal{K} \setminus \{\ell\}$  and satisfaction of the terminal region specification; b) achieving good performance according to the metric  $r(\cdot, \cdot)$ .

Following [3], it can be shown that an infinitesimal expression for a continuously differentiable value function  $V_\ell^u := J_\ell(\cdot, \cdot, u)$ , which is equivalent to (6), is given by

$$\begin{aligned} \text{LE}(V_\ell^u, u) &= \nabla_t V_\ell^u(e_\ell, t) + \nabla_{e_\ell} V_\ell^u(e_\ell, t)^\top (F_\ell(e_\ell) \\ &\quad + G_\ell(e_\ell)u) + \gamma r(e_\ell, u) + L_{k,\ell}(e_\ell) = 0, \end{aligned} \quad (8)$$

which is a partial differential equation. If we let  $t_k = t_\ell$ , then owing to (6) the boundary condition of (8) is

$$V_\ell^u(e_\ell(t_\ell), t_\ell) = \phi(e(t_\ell)). \quad (9)$$

Define the optimal value function as  $V_\ell^*(e_\ell, t) = \min_u J_\ell(e_\ell, t, u)$ , for all  $e_\ell \in \mathbb{R}^n$ ,  $t \in [t_k, t_\ell]$ . If  $V_\ell^*$  is continuously differentiable, by following [3] we can derive the corresponding minimizing controller  $u_\ell^*$  as:

$$u_\ell^*(e_\ell, t) = -\frac{1}{2\gamma} R^{-1} G_\ell(e_\ell)^\top \nabla_{e_\ell} V_\ell^*(e_\ell, t). \quad (10)$$

Combining (10) with  $\text{LE}(V_\ell^*, u_\ell^*) = 0$ , we obtain the Hamilton-Jacobi-Bellman (HJB) equation:

$$\begin{aligned} \nabla_t V_\ell^*(e_\ell, t) + \gamma Q(e_\ell) + L_{k,\ell}(e_\ell) + \nabla_{e_\ell} V_\ell^*(e_\ell, t)^\top F_\ell(e_\ell) \\ - \frac{1}{4\gamma} \nabla_{e_\ell} V_\ell^*(e_\ell, t)^\top G_\ell(e_\ell) R^{-1} G_\ell(e_\ell)^\top \nabla_{e_\ell} V_\ell^*(e_\ell, t) = 0, \\ V_\ell^*(e_\ell, t_\ell) = \phi(e_\ell). \end{aligned} \quad (11)$$

Equation (11) needs to be solved in order to compute (10).

The following theorem shows that if  $\gamma$  is picked sufficiently small and a controllability condition holds, then the optimal policy  $u_\ell^*$  can achieve an optimal timed transition.

**Theorem 1.** Assume that there exists a control law  $u_\ell^c : \mathbb{R}^n \times [t_k, t_\ell]$  such that the closed-loop system  $e_\ell(t)$  under  $u = u_\ell^c$  satisfies: a)  $e_\ell(t_\ell) = 0$ ; b)  $L_{k,\ell}(e_\ell(t)) = 0$ , for all  $t \in [t_k, t_\ell]$ . Then, there exists  $\gamma^* > 0$ , such that if  $\gamma < \gamma^*$ , the closed-loop system  $e_\ell(t)$  under  $u = u_\ell^*$  executes an optimal timed transition, according to Def. 9.

*Proof.* Denote by  $e_\ell^c$  the trajectories of  $e_\ell$  under  $u = u_\ell^c$ , and by  $e_\ell^*$  the trajectories of  $e_\ell$  under  $u = u_\ell^*$ . By definition, it holds that  $L_{k,\ell}(e_\ell^c(t)) = 0$  for all  $t \in [t_k, t_\ell]$ , and  $\phi(e_\ell^c(t_\ell)) = 0$ . Hence,

$$J_\ell(e_\ell(t_k), t_k, u_\ell^c) = \int_{t_k}^{t_\ell} \gamma r(e_\ell^c(\tau), u_\ell^c(e_\ell^c(\tau))) d\tau. \quad (12)$$

By optimality, it follows that

$$0 \leq J_\ell(e_\ell(t_k), t_k, u_\ell^*) \leq J_\ell(e_\ell(t_k), t_k, u_\ell^c). \quad (13)$$

Due to (12),  $\lim_{\gamma \rightarrow 0^+} J_\ell(e_\ell(t_k), t_k, u_\ell^c) = 0$ . As a result, it follows from (13) that  $\lim_{\gamma \rightarrow 0^+} J_\ell(e_\ell(t_k), t_k, u_\ell^*) = 0$ , hence also  $\lim_{\gamma \rightarrow 0^+} L_{k,\ell}(e_\ell^*(t)) = 0$  for all  $t \in [t_k, t_\ell]$ , and  $\lim_{\gamma \rightarrow 0^+} \phi(e_\ell^*(t_\ell)) = 0$ . Hence, for all  $\epsilon^* > 0$  there exists  $\gamma^* > 0$ , such that if  $\gamma < \gamma^*$ , then  $L_{k,\ell}(e_\ell^*(t)) < \epsilon^*$  for all  $t \in [t_k, t_\ell]$  and  $\phi(e_\ell^*(t_\ell)) < \epsilon^*$ . The result follows by the design properties of  $L_{k,\ell}$  and  $\phi$ .  $\blacksquare$

In what follows, we drop the subscript  $\ell$  for ease of exposition.

### B. Policy Iteration

Equation (11), which has to be solved in order to obtain the optimal controller (10), is highly nonlinear. Therefore, its analytic solution is hard to obtain, and we have to resort to approximate solution methods. To do so, we will require a few definitions and assumptions. First, for the cost (6) to be properly defined and for the corresponding value function to be continuously differentiable, we only consider control laws that are *admissible*, as per the following definition.

**Definition 10.** A control law  $u : \mathbb{R}^n \times [t_k, t_\ell] \rightarrow \mathbb{R}^m$  is *admissible* with respect to the cost (6), denoted by  $u \in \mathcal{U}$ , if

- $u$  is continuous over  $\mathbb{R}^n \times [t_k, t_\ell]$  with  $u(0, t) = 0$  for all  $t \in [t_k, t_\ell]$ ; and
- the origin of system (7) is uniformly Lyapunov stable under  $u$ , the trajectories of (7) are bounded for all  $t \in [t_k, t_\ell]$ , and the cost (6) is finite for all  $e_\ell, t_k$ .  $\square$

Next, let  $\mathcal{P}_+$  denote the set of continuously differentiable functions from  $\mathbb{R}^n \times [t_k, t_\ell]$  to  $\mathbb{R}$ . For any function  $V \in \mathcal{P}_+$ , assume that  $V(\cdot, \bar{t})$  is positive-definite for all  $\bar{t} \in [t_k, t_\ell]$ . We need the following assumption for the optimal value function, which is standard in the literature [16], [18].

**Assumption 2.** The optimal value function  $V^*$ , which solves the HJB equation (11), belongs to  $\mathcal{P}_+$ , i.e.,  $V^* \in \mathcal{P}_+$ .  $\square$

Next, we present the Policy Iteration (PI) algorithm for solving the finite-horizon, time-varying HJB equation.

### Policy Iteration

Let  $u_0 \in \mathcal{U}$ . Then, for all  $i \in \mathbb{N}$ , perform the iteration:

- 1) Evaluate the value function  $V^{u_i}$  by solving (8):

$$\nabla_t V^{u_i}(e, t) + \nabla_e V^{u_i}(e, t)^\top (F(e) + G(e)u_i) + \gamma r(e, u_i) + L(e) = 0, \quad \forall t \in [t_k, t_\ell], \quad (14)$$

with  $V^{u_i}(e(t_\ell), t_\ell) = \phi(e(t_\ell))$ .

- 2) Choose the next control law  $u_{i+1}$  as

$$u_{i+1}(e, t) = -\frac{1}{2\gamma} R^{-1} G(e)^\top \nabla_e V^{u_i}(e, t). \quad (15)$$

The following lemma is crucial to establishing convergence of the PI algorithm.

**Lemma 1.** Consider the sequence of control laws  $\{u_i\}_{i \in \mathbb{N}}$  and continuously differentiable value functions  $\{V^{u_i}\}_{i \in \mathbb{N}}$  generated by the PI algorithm through equations (14)-(15). Let  $u_i$  be admissible, for some  $i \in \mathbb{N}$ . Then:

- 1)  $u_{i+1}$  is admissible.
- 2)  $V^*(e, t) \leq V^{u_{i+1}}(e, t) \leq V^{u_i}(e, t)$ ,  $\forall (e, t) \in \mathbb{R}^n \times [t_k, t_\ell]$ .

*Proof.* For the first part, notice that  $G$  is continuous by assumption, and  $\nabla_e V^{u_i}$  is also continuous since  $V^{u_i}$  is assumed to be continuously differentiable. Therefore,  $u_{i+1}$  is continuous. In addition, owing to the admissibility of  $u_i$ , it holds that  $V^{u_i}(0, t) = 0$ , for all  $t \in [t_k, t_\ell]$ . Hence, since  $V^{u_i}$  is non-negative, it follows that  $\nabla_e V^{u_i}(0, t) = 0$  because  $V^{u_i}$  attains a minimum at  $(0, t)$ , which yields  $u_{i+1}(0, t) = 0$ , for all  $t \in [t_k, t_\ell]$ . Next, over the trajectories of (3) and under the control law  $u = u_{i+1}$ , one has

$$\begin{aligned} \dot{V}^{u_i} &= \nabla_t V^{u_i} + (\nabla_e V^{u_i})^\top (F + Gu_i) + (\nabla_e V^{u_i})^\top G(u_{i+1} - u_i) \\ &= -\gamma Q(e) - \gamma S(u_i) - L(e) - 2\gamma u_{i+1}^\top R(u_{i+1} - u_i) \\ &= -\gamma Q(e) - L(e) - \gamma S(u_{i+1}) - \gamma S(u_{i+1} - u_i) \leq 0, \end{aligned} \quad (16)$$

where we used (14) and (15). Moreover, for any  $(e, t) \in \mathbb{R}^n \times [t_k, t_\ell]$  it holds that  $V^{u_i}(e, t) \geq V^*(e, t)$ . Since  $V^* \in \mathcal{P}_+$ , for any fixed  $\bar{t} \in [t_k, t_\ell]$  there exists a class  $\mathcal{K}_\infty$  function  $a_1^{\bar{t}}$  so that  $a_1^{\bar{t}}(\|e\|) \leq V^*(e, \bar{t})$ . Hence we deduce the existence of a class  $\mathcal{K}_\infty$  function  $a_1$  such that  $V^*(e, t) \geq a_1(\|e\|) := \min\{a_1^{\bar{t}}(\|e\|), \bar{t} \in [t_k, t_\ell]\}$ . Combining this result with (16) we obtain  $a_1(\|e\|) \leq V^*(e, t) \leq V^{u_i}(e, t) \leq V^{u_i}(e_k, t_k)$ , hence  $\|e\| \leq a_1^{-1}(V^{u_i}(e_k, t_k))$ , proving boundedness of the system's trajectories as well as the finiteness of (6). Finally,

it can be seen that the function  $W(e) := \max\{V^{u_i}(e, \bar{t}), \bar{t} \in [t_k, t_\ell]\}$  satisfies  $W(0) = 0$  owing to  $F(0) = 0$ ,  $L(0) = 0$  and  $u(0, t) = 0$  for all  $t \in [t_k, t_\ell]$ , and  $W(e) > 0$  for any  $x \neq 0$  owing to the positivity of the running and the terminal cost in (6), and  $V^{u_i}(e, t) \leq W(e)$  for all  $t \in [t_k, t_\ell]$ . Hence, it follows from [24] that the origin of (7) is uniformly Lyapunov stable under the control  $u_{i+1}$ .

For item 2), integration of (16) over  $t \in [t_k, t_\ell]$  yields:

$$\begin{aligned} V^{u_i}(e(t_\ell), t_\ell) - V^{u_i}(e_k, t_k) &= - \int_{t_k}^{t_\ell} \left( \gamma Q(e) \right. \\ &\quad \left. + \gamma S(u_{i+1}) + L(e) \right) d\tau - \int_{t_k}^{t_\ell} \gamma S(u_{i+1} - u_i) d\tau. \end{aligned} \quad (17)$$

Owing to (9), we have  $V^{u_i}(e(t_\ell), t_\ell) = V^{u_{i+1}}(e(t_\ell), t_\ell) = \phi(e(t_\ell))$ . Therefore, (17) is equivalent to:

$$V^{u_i}(e_k, t_k) = V^{u_{i+1}}(e_k, t_k) + \int_{t_k}^{t_\ell} \gamma S(u_{i+1} - u_i) d\tau$$

Hence,  $V^{u_{i+1}}(e, t) \leq V^{u_i}(e, t)$ , for all  $(e, t) \in \mathbb{R}^n \times [t_k, t_\ell]$ , while the inequality  $V^*(e, t) \leq V^{u_i}(e, t)$  holds by optimality.  $\blacksquare$

**Theorem 2.** Let  $u_0 \in \mathcal{U}$ . Then, the PI algorithm described through equations (14)-(15) guarantees that  $\lim_{i \rightarrow \infty} V^{u_i} = V^*$  and  $\lim_{i \rightarrow \infty} u_i = u^*$ . The convergence is uniform on any compact subset of  $\mathbb{R}^n \times [t_k, t_\ell]$ .

*Proof.* Given the monotonicity results of Lemma 1, the proof follows similar steps with [25] and is thus omitted.  $\blacksquare$

### C. Approximate Solution to the Time-Varying HJB Equation

The PI algorithm requires knowledge of the system's dynamics functions  $F$ ,  $G$ . Towards implementing a model-free version of PI, we rewrite the system error dynamics as

$$\dot{e} = F(e) + G(e)u_i(e, t) + G(e)v_i(e, t), \quad t \geq 0, \quad (18)$$

where  $v_i = u - u_i$ ,  $i \in \mathbb{N}$ , and  $u_i$  is as defined in (15). Taking the total time derivative of the value function  $V^{u_i}$ ,  $i \in \mathbb{N}$ , along the trajectories of (18), and using (14)-(15), we obtain

$$\begin{aligned} \dot{V}^{u_i} &= \nabla_t V^{u_i} + (\nabla_e V^{u_i})^\top (F(e) + G(e)u_i(e, t) + G(e)v_i(e, t)) \\ &= -\gamma Q(e) - \gamma S(u) - L(e) - 2\gamma u_{i+1}^\top (e, t)^\top R v_i(e, t). \end{aligned} \quad (19)$$

Integrating (19) over any time interval  $[t, t+T] \subseteq [t_k, t_\ell]$ , with  $T > 0$  and for all  $t \in [t_k, t_l - T]$ , we derive

$$\begin{aligned} V^{u_i}(e(t+T), t+T) - V^{u_i}(e(t), t) &= - \int_t^{t+T} \left( \gamma Q(e) \right. \\ &\quad \left. + L(e) + \gamma S(u_i(e, \tau)) + 2\gamma u_{i+1}^\top(e, \tau)^\top R v_i(e, \tau) \right) d\tau, \end{aligned} \quad (20)$$

$$V^{u_i}(e(t_\ell), t_\ell) = \phi(e(t_\ell)). \quad (21)$$

Notice that (20)-(21) is a model-free version of (14), as it is independent of the functions  $F$ ,  $G$ . However, owing to the infinite dimensionality of this equation, we need to resort to approximation theory in order to solve it with respect to  $u_{i+1}$  and  $V^{u_i}$ . Particularly, we can use the Weierstrass approximation theorem [3] and deduce that  $V^{u_i}$ ,  $u_i$  can be

uniformly approximated on a compact set  $\Omega \times [t_k, t_\ell] =: D$ , with  $\Omega \subset \mathbb{R}^n$ . Then, we can express  $V^{u_i}$ ,  $u_{i+1}$ ,  $\forall i \in \mathbb{N}$ , as

$$V^{u_i}(e, t) = (w_i^v)^T \psi^v(e, t) + \phi(e) + \epsilon_i^v(e, t), \quad (22a)$$

$$u_{i+1}(e, t) = (w_i^u)^T \psi^u(e, t) + \epsilon_i^u(e, t), \quad (22b)$$

where  $w_i^v \in \mathbb{R}^{N_v}$ ,  $w_i^u \in \mathbb{R}^{N_u \times m}$  are weights,  $\psi^v : \mathbb{R}^n \times [t_k, t_\ell] \rightarrow \mathbb{R}^{N_v}$ ,  $\psi^u : \mathbb{R}^n \times [t_k, t_\ell] \rightarrow \mathbb{R}^{N_u}$  are basis functions and  $\epsilon_i^v : \mathbb{R}^n \times [t_k, t_\ell] \rightarrow \mathbb{R}$ ,  $\epsilon_i^u : \mathbb{R}^n \times [t_k, t_\ell] \rightarrow \mathbb{R}^m$  are the approximation errors. The approximation errors  $\epsilon_i^v$ ,  $\epsilon_i^u$  converge to zero, uniformly on  $D$ , as  $N_v, N_u \rightarrow \infty$ .

As  $w_i^v$  and  $w_i^u$  in (22) are unknown, we construct a critic and an actor neural network to approximate  $V^{u_i}$ ,  $u_{i+1}$  as

$$\hat{V}^{u_i}(e, t) := (\hat{w}_i^v)^T \psi^v(e, t) + \phi(e), \quad (23a)$$

$$\hat{u}_{i+1}(e, t) := (\hat{w}_i^u)^T \psi^u(e, t), \quad (23b)$$

where  $\hat{w}_i^v \in \mathbb{R}^{N_v}$ ,  $\hat{w}_i^u \in \mathbb{R}^{N_u \times m}$  are the critic and the actor weights respectively, and  $i \in \mathbb{N}$ . Notice that a bias term has been introduced for the approximation of the value function in (22)-(23). Its purpose is to impose the boundary condition (21) to hold irrespectively of how the critic weights  $\hat{w}_i^u$  are chosen, as long as the basis functions are appropriate.

**Corollary 1.** *Let  $\psi^v(0, t) = 0$  for all  $t \in [t_k, t_\ell]$ , and  $\psi^v(e, t_\ell) = 0$  for all  $e \in \mathbb{R}^n$ . Then, for all  $i \in \mathbb{N}$ , it holds that:  $\hat{V}^{u_i}(e, t_\ell) = \phi(e)$ , for all  $e \in \mathbb{R}^n$ , and  $\hat{V}^{u_i}(0, t) = 0$ , for all  $t \in [t_k, t_\ell]$ .*

*Proof.* The proof follows by direct substitution in (23). ■

Consider now a number of time instances  $\tau_j$ ,  $j \in \{0, \dots, N\} =: \mathcal{N}$  such that  $t_k = \tau_0 < \tau_1 < \dots < \tau_N = t_\ell$ . Using the approximation (23), the left hand side of (20) for  $t = \tau_j$  and  $t + T = \tau_{j+1}$ ,  $j \in \mathcal{N} \setminus \{N\}$ , is approximated as:

$$\hat{V}^{u_i}(e(\tau_{j+1}), \tau_{j+1}) - \hat{V}^{u_i}(e(\tau_j), \tau_j) = \phi(e(\tau_{j+1})) - \phi(e(\tau_j)) + (\hat{w}_i^v)^T (\psi^v(e(\tau_{j+1}), \tau_{j+1}) - \psi^v(e(\tau_j), \tau_j)). \quad (24)$$

In addition, the term  $2\gamma u_{i+1}(e, \tau)^T R v_i(e, \tau)$  at the right hand side of (20) can be approximated using the actor as:

$$2\gamma \hat{u}_{i+1}(e, \tau)^T R \hat{v}_i(e, \tau) = 2\gamma \psi^u(e, \tau)^T \hat{w}_i^u R \hat{v}_i(e, \tau) = 2\gamma \left( (\hat{v}_i(e, \tau)^T R) \otimes \psi^u(e, \tau)^T \right) \text{vec}(\hat{w}_i^u) \quad (25)$$

where  $\hat{v}_i = u - \hat{u}_i$ . Hence, the residual error created by approximating equation (20) through (24)-(25) is

$$e_{j,i} := \hat{V}^{u_i}(e(\tau_{j+1}), \tau_{j+1}) - \hat{V}^{u_i}(e(\tau_j), \tau_j) + \int_{\tau_j}^{\tau_{j+1}} \left( \gamma Q(e) + L(e) + \gamma S(\hat{u}_i(e, \tau)) + 2\gamma \hat{u}_{i+1}(e, \tau)^T R \hat{v}_i(e, \tau) \right) d\tau,$$

which can be written as:

$$e_{j,i} = \Theta_{j,i} \hat{W}_i + \Psi_{j,i}, \quad (26)$$

where  $\Theta_{j,i} := [\Theta_{j,i}^v \ \Theta_{j,i}^u]$ ,  $\hat{W}_i := [(\hat{w}_i^v)^T \ \text{vec}(\hat{w}_i^u)^T]^T$ , and

$$\Theta_{j,i}^v := \left( \psi^v(e(\tau_{j+1}), \tau_{j+1}) - \psi^v(e(\tau_j), \tau_j) \right)^T, \\ \Theta_{j,i}^u := \int_{\tau_j}^{\tau_{j+1}} 2\gamma \left( (\hat{v}_i(e, \tau)^T R) \otimes \psi^u(e, \tau)^T \right) d\tau,$$

---

### Algorithm 1 Model-Free PI

---

- 1: Employ an arbitrary behavioral policy  $u_b$  to the system (3), and collect input-state data online.
  - 2: Let  $u_0 \in \mathcal{U}$  be admissible, select  $\epsilon > 0$  and set  $i = 0$ .
  - 3: **repeat**
  - 4:   Solve for  $\hat{w}_i^v$  and  $\hat{w}_i^u$  from equation (27) and  $i = i + 1$ .
  - 5: **until**  $\|\hat{w}_i^v - \hat{w}_{i-1}^v\| < \epsilon$
  - 6: Switch from  $u_b$  to the learnt control policy  $\hat{u}_i$ .
- 

$$\Psi_j := \phi(e(\tau_{j+1})) - \phi(e(\tau_j)) + \int_{\tau_j}^{\tau_{j+1}} \left( \gamma Q(e) + L(e) + \gamma S(\hat{u}_i(e, \tau)) \right) d\tau.$$

If enough data is obtained along the system's trajectories, we can find  $\hat{W}_i$  by least squares to minimize the error (26). To that end, we impose a standard assumption [17], [18].

**Assumption 3.** There exist  $\delta > 0$  and  $l_0 \in \mathcal{N}$ , so that for all  $l \geq l_0$  it holds that  $\sum_{j=0}^l \Theta_{j,i}^T \Theta_{j,i} > l\delta I_{N_v + mN_u}$ . □

Given Assumption 3, the least squares solution to (26) is:

$$\hat{W}_i = - \left( \sum_{j=0}^l \Theta_{j,i}^T \Theta_{j,i} \right)^{-1} \left( \sum_{j=0}^l \Theta_{j,i}^T \Psi_{j,i} \right). \quad (27)$$

As a result, we can obtain the model-free version of PI, as shown in Algorithm 1. Its convergence is shown next.

**Theorem 3.** *Let Assumption 3 hold. Then, for all  $\epsilon > 0$  there exist constants  $N_v^m$ ,  $N_u^m$ ,  $i^* \in \mathbb{N}$ , such that if  $N_v \geq N_v^m$  and  $N_u \geq N_u^m$ , then for all  $(e, t) \in D$ ,  $i \geq i^*$ , it holds that*

$$\|\hat{V}^{u_i}(e, t) - V^*(e, t)\| \leq \epsilon, \quad \|\hat{u}_{i+1}(e, t) - u^*(e, t)\| \leq \epsilon.$$

*Proof.* For  $i \in \mathbb{N}$ , let  $\tilde{V}^{u_i}$  be the value function of  $\hat{u}_i$ , where  $\hat{u}_0 = u_0$ , so that  $\text{LE}(\tilde{V}^{u_i}, \hat{u}_i) = 0$ ,  $\tilde{V}^{u_i}(0, t) = 0$ , for all  $t \in [t_k, t_\ell]$ , and  $\tilde{V}^{u_i}(e, t_\ell) = \phi(e)$ , for all  $e \in \Omega$ . Let also  $\tilde{u}_{i+1}(e, t) = -\frac{1}{2\gamma} R^{-1} G(e)^T \nabla_e \tilde{V}^{u_i}(e, t)$ ,  $\forall (e, t) \in D$ . Then, using the integral form of expression for value functions (20) over the trajectories of (3), it follows for all  $j \in \mathcal{N}$  that

$$\tilde{V}^{u_i}(e(\tau_{j+1}), \tau_{j+1}) - \tilde{V}^{u_i}(e(\tau_j), \tau_j) = - \int_{\tau_j}^{\tau_{j+1}} \left( \gamma Q(e) + L(e) + \gamma S(u_i(e, \tau)) + 2\gamma \tilde{u}_{i+1}^T(e, \tau) R \hat{v}_i(e, \tau) \right) d\tau \quad (28)$$

The value function  $\tilde{V}^{u_i}$  and the controller  $\tilde{u}_{i+1}$  can be uniformly approximated on  $D$ . Hence, there exist  $\tilde{w}_i^v \in \mathbb{R}^{N_v}$ ,  $\tilde{w}_i^u \in \mathbb{R}^{N_u \times m}$  such that  $\tilde{V}^{u_i}(e, t) = (\tilde{w}_i^v)^T \psi^v(e, t) + \phi(e) + \tilde{\epsilon}_i^v(e, t)$  and  $\tilde{u}_{i+1}(e, t) = (\tilde{w}_i^u)^T \psi^u(e, t) + \tilde{\epsilon}_i^u(e, t)$ . The approximation errors  $\tilde{\epsilon}_i^v : \mathbb{R}^n \times [t_k, t_\ell] \rightarrow \mathbb{R}$ ,  $\tilde{\epsilon}_i^u : \mathbb{R}^n \times [t_k, t_\ell] \rightarrow \mathbb{R}^m$  vanish uniformly on  $D$  as  $N_v, N_u \rightarrow \infty$ . Substituting these expressions in (28), we have:

$$0 = \Theta_{j,i} \tilde{W}_i + \Psi_{j,i} + E_{j,i}, \quad \forall j \in \mathcal{N}, \quad i \in \mathbb{N}, \quad (29)$$

where  $\tilde{W}_i = [\tilde{w}_i^v{}^T \ \text{vec}(\tilde{w}_i^u)^T]^T$  and  $E_{j,i} = \tilde{\epsilon}_i^v(e(\tau_{j+1}), \tau_{j+1}) - \tilde{\epsilon}_i^v(e(\tau_j), \tau_j) + \int_{\tau_j}^{\tau_{j+1}} 2\gamma \tilde{\epsilon}_i^u(e, \tau)^T R \hat{v}_i(e, \tau) d\tau$ . Deducting (29) from (26) yields  $e_{j,i} = \Theta_{j,i} \tilde{W}_i - E_{j,i}$ , where  $\tilde{W}_i = \hat{W}_i - \tilde{W}_i$ . Due to Assumption 3, the least squares estimate  $\hat{W}_i$

from equation (27) is well-defined and bounded. As a result, the least squares algorithm yields  $\sum_{j=0}^l e_{j,i}^2 \leq \sum_{j=0}^l E_{j,i}^2$ . In addition, since  $\sum_{j=0}^l \bar{W}_i^T \Theta_{j,i}^T \Theta_{j,i} \bar{W}_i = \sum_{j=0}^l (e_{j,i} - E_{j,i})^2$ , using the condition of Assumption 3 we derive  $\|\bar{W}_i\|^2 \leq \frac{4}{\delta} \max_{0 \leq j \leq l} E_{j,i}^2$ . Due to its dependence on  $\bar{\epsilon}_i^v, \bar{\epsilon}_i^u$ , it follows that as  $N_v, N_u \rightarrow \infty$ ,  $E_{j,i}^2$  converges uniformly to zero on  $D$  for all  $i \in \mathbb{N}, j \in \mathcal{N}$ , hence  $\bar{W}_i$  also converges uniformly to zero on  $D$ . Thus, for all  $\epsilon > 0$  there exist  $N_v^*, N_u^* > 0$ , such that if  $N_v \geq N_v^*, N_u \geq N_u^*$  then  $\forall (e, t) \in D$  it holds

$$\begin{aligned} |\hat{V}^{u_i}(e, t) - \tilde{V}^{u_i}(e, t)| &= |(\hat{w}_i^v - \tilde{w}_i^v)| |\psi_i^v(e, t)| \quad (30) \\ &+ |\bar{\epsilon}_i^v(e, t)| \leq \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon, \end{aligned}$$

$$\begin{aligned} |\hat{u}_i(e, t) - \tilde{u}_i(e, t)| &= |\hat{w}_i^u - \tilde{w}_i^u| |\psi_i^u(e, t)| \quad (31) \\ &+ |\bar{\epsilon}_i^u(e, t)| \leq \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon. \end{aligned}$$

Finally, we use an induction to derive the final result.

1) For  $i = 0$ , we have  $\tilde{V}^{u_0} = V^{u_0}$  and  $\tilde{u}_1 = u_1$ . Hence, due to the uniform convergence (30)-(31), it follows that  $\lim_{N_v, N_u \rightarrow \infty} \hat{V}^{u_0}(e, t) = V^{u_0}(e, t)$  and  $\lim_{N_v, N_u \rightarrow \infty} \hat{u}_1(e, t) = u_1(e, t)$ , uniformly on  $D$ .

2) Suppose that  $\lim_{N_v, N_u \rightarrow \infty} \hat{V}^{u_{i-1}}(e, t) = V^{u_{i-1}}(e, t)$  and  $\lim_{N_v, N_u \rightarrow \infty} \hat{u}_i(e, t) = u_i(e, t)$ , uniformly on  $D$ , for some  $i \in \mathbb{N}$ . Then, exploiting (20)-(21), we have

$$\begin{aligned} & \left| \tilde{V}^{u_i}(e(t_\ell), t_\ell) - \tilde{V}^{u_i}(e(t), t) - (V^{u_i}(e(t_\ell), t_\ell) \right. \\ & \quad \left. - V^{u_i}(e(t), t)) \right| = \left| V^{u_i}(e(t), t) - \tilde{V}^{u_i}(e(t), t) \right| = \\ & \left| \int_t^{t_\ell} \gamma \left( S(u_i(e, \tau)) + 2u_{i+1}^T(e, \tau) R v_i(e, \tau) \right) d\tau \right. \\ & \quad \left. - \int_t^{t_\ell} \gamma \left( S(\hat{u}_i(e, \tau)) + 2\hat{u}_{i+1}^T(e, \tau) R \hat{v}_i(e, \tau) \right) d\tau \right| \\ & \leq \left| \int_t^{t_\ell} \gamma \left( S(u_i(e, \tau)) - S(\hat{u}_i(e, \tau)) \right) d\tau \right| \\ & + \left| \int_t^{t_\ell} 2\gamma \left( (u_{i+1}(e, \tau) - \hat{u}_{i+1}(e, \tau))^T R \hat{v}_i(e, \tau) \right) d\tau \right| \\ & + \left| \int_t^{t_\ell} \left( 2\gamma u_{i+1}^T(e, \tau) R (u_i(e, \tau) - \hat{u}_i(e, \tau)) \right) d\tau \right| \quad (32) \end{aligned}$$

However, due to the inductive assumptions we know that

$$\begin{aligned} \lim_{N_u, N_v \rightarrow \infty} \int_t^{t_\ell} \gamma \left( S(u_i(e, \tau)) - S(\hat{u}_i(e, \tau)) \right) d\tau &= 0, \\ \lim_{N_u, N_v \rightarrow \infty} \int_t^{t_\ell} \left( 2\gamma u_{i+1}^T(e, \tau) R (u_i(e, \tau) - \hat{u}_i(e, \tau)) \right) d\tau &= 0, \end{aligned}$$

uniformly on  $D$ . In addition, due to Assumption 3, it holds that  $\lim_{N_u, N_v \rightarrow \infty} \tilde{u}_{i+1}(e, t) = u_{i+1}(e, t)$ . Therefore, due to the three limits, (32) yields  $\lim_{N_u, N_v \rightarrow \infty} \tilde{V}^{u_i}(e, t) = V^{u_i}(e, t)$ . Hence, since  $|V^{u_i}(e, t) - \hat{V}^{u_i}(e, t)| \leq |V^{u_i}(e, t) - \tilde{V}^{u_i}(e, t)| + |\tilde{V}^{u_i}(e, t) - \hat{V}^{u_i}(e, t)|$ , we can use the inductive assumption to conclude that,  $\forall \epsilon > 0$ , there exist  $N_v^{**}, N_u^{**} > 0$  such that if  $N_v \geq N_v^{**}, N_u \geq N_u^{**}$  then  $\forall (e, t) \in D$ :

$$|V^{u_i}(e, t) - \hat{V}^{u_i}(e, t)| \leq \epsilon, \quad (33)$$

which concludes the induction. The result follows from (33) and Theorem 2, by using the triangular inequality. ■

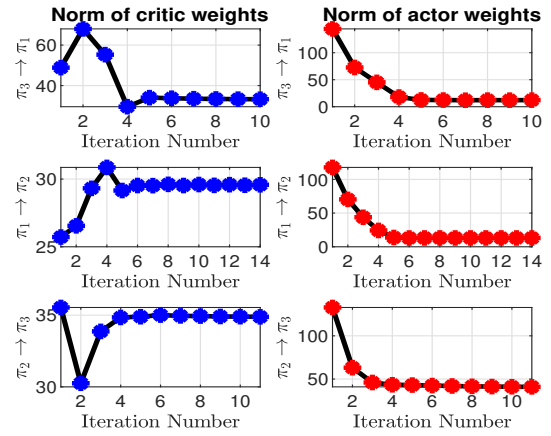


Fig. 1. Evolution of the Frobenius norms of the actor and the critic weights, as derived by Alg. 1.

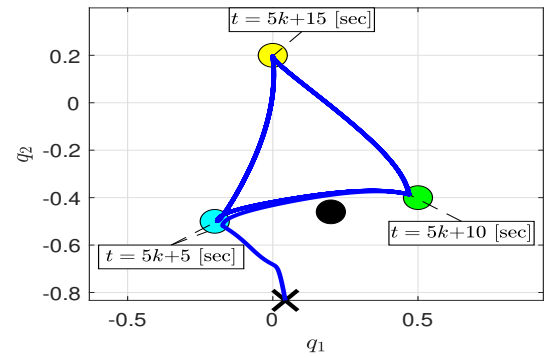


Fig. 2. Evolution of  $q$  after employing the policy given by Alg. 1.

*Remark 1.* Due to Theorems 1 and 3, if  $N_v, N_u$  are large enough and  $\gamma$  is small, the closed-loop system eventually guarantees an  $\epsilon$ -optimal safe timed transition between  $\pi_k$  and  $\pi_\ell$  as per Def. 9. Since the aforementioned results apply for the transition among any pair of regions, we conclude that the closed-loop system eventually satisfies  $\varphi$  safely. □

## VI. SIMULATIONS

We consider a two-link manipulator [26],  $M(q)\ddot{q} + V_m(q, \dot{q})\dot{q} + F_d\dot{q} + F_s(q) = u$ , where  $q = [q_1 \ q_2]^T$  and  $\dot{q} = [\dot{q}_1 \ \dot{q}_2]^T$  are the angular positions (in rad) and the angular velocities (in rad/s), respectively. The matrices  $M(q) \in \mathbb{R}^{2 \times 2}$  and  $V_m(q, \dot{q}) \in \mathbb{R}^{2 \times 2}$  are the inertia and the centripetal-Coriolis matrices, while  $F_d\dot{q}$  and  $F_s(q)$  model the dynamic and static friction, respectively [26]. We consider three regions of interest  $\Pi = \{\pi_1, \dots, \pi_3\}$  centered at  $c_1 = [-0.2, -0.2, 0, 0]^T$ ,  $c_2 = [0.5, -0.4, 0, 0]^T$ ,  $c_3 = [0, 0.2, 0, 0]^T$ , and a joint-state obstacle centered at  $o_4 = [0.2, -0.46]^T$ , all with radius 0.05. Further, we consider  $\mathcal{AP} := \{‘1’, ‘2’, ‘3’\}$  and  $\mathcal{L}(\pi_i) = \{‘i’\}$ ,  $i \in \{1, \dots, 3\}$ .

We impose a timed temporal logic task dictated by the formula  $\varphi = \square \diamond_{[0,5]} ‘i’, i \in \{1, 2, 3\}$ , implying periodic visit to regions  $\pi_1, \pi_2, \pi_3$  every 5 seconds; we also require avoidance of the obstacle, for which we compute  $L$  using  $o_4$ . By setting  $\gamma(\cdot) = 5$  for all transitions in  $\Pi \times \Pi$  in (5), and following the methodology of Section IV, we

obtain the repetitive timed path  $\mathbf{p} = [(\pi_1, 5k + 5)(\pi_2, 5k + 10)(\pi_3, 5k + 15)]^\omega$  for  $k \in \{0, 1, \dots\}$ . We perform Alg. 1 by employing a sinusoidal behavioral policy  $u_b$  for 150 seconds, and then executing the model-free PI by solving Eq. (27) iteratively. The evolution of the critic-actor weight norms during the execution of Alg. 1 are illustrated in Fig. 1, showing their convergence. After the passage of the 150 seconds, the policy derived by Alg. 1 substitutes the behavioral policy, and the resulting closed-loop trajectories for  $t \geq 150$  [sec] can be seen in Fig. 2. It can be verified that the closed-loop system executes successfully the timed path, leading to the eventual satisfaction of  $\varphi$ . For all three repetitive optimal control problems, we chose  $R=0.5I_2$ ,  $\phi(e)=Q(e)=e^T(\text{diag}[20 \ 20 \ 10 \ 10])e$ , and  $\gamma=0.1$ .

## VII. CONCLUSION

We develop a two-layered algorithm for the planning and control of unknown systems with timed temporal logic tasks. We design a novel data-driven control protocol that learns how to execute optimal timed transition between regions of the state-space, which guarantees the eventual satisfaction of the task. Future efforts will be devoted towards addressing continuous-time temporal tasks under the same framework.

## REFERENCES

- [1] C. Baier and J.-P. Katoen, *Principles of model checking*. MIT press, 2008.
- [2] C. K. Verginis and D. V. Dimarogonas, "Timed abstractions for distributed cooperative manipulation," *Autonomous Robots*, vol. 42, no. 4, pp. 781–799, 2018.
- [3] F. L. Lewis, D. Vrabie, and V. L. Syrmos, *Optimal control*. John Wiley & Sons, 2012.
- [4] S. Karaman and E. Frazzoli, "Vehicle routing problem with metric temporal logic specifications," in *2008 47th IEEE Conference on Decision and Control*. IEEE, 2008, pp. 3953–3958.
- [5] C. C. Constantinou and S. G. Loizou, "Automatic controller synthesis of motion-tasks with real-time objectives," in *2018 IEEE Conference on Decision and Control (CDC)*. IEEE, 2018, pp. 403–408.
- [6] Z. Lin and J. S. Baras, "Metric interval temporal logic based reinforcement learning with runtime monitoring and self-correction," in *2020 American Control Conference (ACC)*. IEEE, 2020, pp. 5400–5406.
- [7] C. N. Mavridis, C. Vrohidis, J. S. Baras, and K. J. Kyriakopoulos, "Robot navigation under mtl constraints using time-dependent vector field based control," in *2019 IEEE 58th Conference on Decision and Control (CDC)*. IEEE, 2019, pp. 232–237.
- [8] P. Varnai and D. V. Dimarogonas, "Prescribed performance control guided policy improvement for satisfying signal temporal logic tasks," in *2019 American Control Conference (ACC)*. IEEE, 2019, pp. 286–291.
- [9] A. Nikou, D. Boskos, J. Tumova, and D. V. Dimarogonas, "On the timed temporal logic planning of coupled multi-agent systems," *Automatica*, vol. 97, pp. 339–345, 2018.
- [10] A. Nikou, S. Heshmati-Alamdari, C. K. Verginis, and D. V. Dimarogonas, "Decentralized abstractions and timed constrained planning of a general class of coupled multi-agent systems," in *2017 IEEE 56th Annual Conference on Decision and Control (CDC)*. IEEE, 2017, pp. 990–995.
- [11] C. K. Verginis, C. Vrohidis, C. P. Bechlioulis, K. J. Kyriakopoulos, and D. V. Dimarogonas, "Reconfigurable motion planning and control in obstacle cluttered environments under timed temporal tasks," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 951–957.
- [12] L. Lindemann and D. V. Dimarogonas, "Control barrier functions for signal temporal logic tasks," *IEEE Control Systems Letters*, vol. 3, no. 1, pp. 96–101, 2018.
- [13] W. Xiao, C. A. Belta, and C. G. Cassandras, "High order control lyapunov-barrier functions for temporal logic specifications," *arXiv preprint arXiv:2102.06787*, 2021.

- [14] C. Sun and K. G. Vamvoudakis, "Continuous-time safe learning with temporal logic constraints in adversarial environments," in *2020 American Control Conference (ACC)*. IEEE, 2020, pp. 4786–4791.
- [15] D. Muniraj, K. G. Vamvoudakis, and M. Farhood, "Enforcing signal temporal logic specifications in multi-agent adversarial environments: A deep q-learning approach," in *2018 IEEE Conference on Decision and Control (CDC)*. IEEE, 2018, pp. 4141–4146.
- [16] K. G. Vamvoudakis and F. L. Lewis, "Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem," *Automatica*, vol. 46, no. 5, pp. 878–888, 2010.
- [17] Y. Jiang and Z.-P. Jiang, "Robust adaptive dynamic programming and feedback stabilization of nonlinear systems," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 25, no. 5, pp. 882–893, 2014.
- [18] W. Gao and Z.-P. Jiang, "Learning-based adaptive optimal tracking control of strict-feedback nonlinear systems," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 6, pp. 2614–2624, 2017.
- [19] R. Alur and D. L. Dill, "A theory of timed automata," *Theoretical Computer Science*, vol. 126, no. 2, pp. 183–235, 1994.
- [20] P. Bouyer, N. Markey, J. Ouaknine, and J. Worrell, "The cost of punctuality," in *22nd Annual IEEE Symposium on Logic in Computer Science (LICS 2007)*. IEEE, 2007, pp. 109–120.
- [21] C.-I. Vasile, D. Aksaray, and C. Belta, "Time window temporal logic," *Theoretical Computer Science*, vol. 691, pp. 27–54, 2017.
- [22] D. D'Souza and P. Prabhakar, "On the expressiveness of mtl in the pointwise and continuous semantics," *International Journal on Software Tools for Technology Transfer*, vol. 9, no. 1, pp. 1–4, 2007.
- [23] J. Ouaknine and J. Worrell, "On the decidability of metric temporal logic," in *20th Annual IEEE Symposium on Logic in Computer Science (LICS'05)*. IEEE, 2005, pp. 188–197.
- [24] H. K. Khalil, *Nonlinear systems*. Prentice hall Upper Saddle River, NJ, 2002, vol. 3.
- [25] M. Abu-Khalaf and F. L. Lewis, "Nearly optimal control laws for nonlinear systems with saturating actuators using a neural network hjb approach," *Automatica*, vol. 41, no. 5, pp. 779–791, 2005.
- [26] R. Kamalapurkar, H. Dinh, S. Bhasin, and W. E. Dixon, "Approximate optimal trajectory tracking for continuous-time nonlinear systems," *Automatica*, vol. 51, pp. 40–48, 2015.

## APPENDIX

### Design of the penalty term $L_{k,\ell}$

Let  $k, \ell \in \mathcal{K}$  with  $k \neq \ell$ . We first unify the notation for unsafe zones and regions for simplicity. Let  $\tilde{\Pi} := \{\tilde{\pi}_1, \dots, \tilde{\pi}_K, \tilde{\pi}_{K+1}, \dots, \tilde{\pi}_{K+K_o}\}$ ,  $\tilde{K} := \{1, \dots, K + K_o\}$ , and  $\tilde{\pi}_i := \tilde{B}(\tilde{c}_i, \tilde{\rho}_i)$ , with  $\tilde{c}_i = c_i$ ,  $\tilde{\rho}_i = \rho_i$ , for  $i \in \mathcal{K}$ , and  $\tilde{c}_i = c_{o_i-K}$ ,  $\tilde{\rho}_i = \rho_{o_i-K}$ , for  $i \in \{K + 1, \dots, K + K_o\}$ .

Select  $\bar{\rho}$  such that  $\|\tilde{c}_i - \tilde{c}_j\| \geq \tilde{\rho}_i + \tilde{\rho}_j + \bar{\rho}$ , for all  $i, j \in \tilde{\mathcal{K}} \setminus \{k, \ell\}$ , with  $i \neq j'$ . Note that such a  $\bar{\rho}$  exists, since the regions in  $\Pi \cup \mathcal{O}$  are pairwise disjoint. Let now  $L_i : [-\tilde{\rho}_i^2, \infty) \rightarrow \mathbb{R}_{\geq 0}$  be twice differentiable *non-increasing* functions satisfying  $L_i(x) = \bar{L}$  for all  $x \in [-\tilde{\rho}_i^2, 0]$  for some  $\bar{L} > 0$ , and  $L_i(x) = 0$ , for all  $x \geq \bar{\rho}^2$  and all  $i \in \tilde{\mathcal{K}} \setminus \{k, \ell\}$ . Then it can be proved that  $L_i(\|x - \tilde{c}_i\|^2 - \tilde{\rho}_i^2) \neq 0$  for some  $i \in \tilde{\mathcal{K}} \setminus \{k, \ell\}$  implies  $L_j(\|x - \tilde{c}_j\|^2 - \tilde{\rho}_j^2) = 0$  for all  $j \in \tilde{\mathcal{K}} \setminus \{i, k, \ell\}$ . Indeed,  $L_i(\|x - \tilde{c}_i\|^2 - \tilde{\rho}_i^2) \neq 0$  implies that  $\|x - \tilde{c}_i\|^2 - \tilde{\rho}_i^2 < \bar{\rho}^2$ , i.e.,  $\|x - \tilde{c}_i\| \leq \bar{\rho} + \tilde{\rho}_i < \|\tilde{c}_i - \tilde{c}_j\|$  for all  $j \in \tilde{\mathcal{K}} \setminus \{i, k, \ell\}$ . Therefore,  $\|x - \tilde{c}_j\| \geq \|\tilde{c}_i - \tilde{c}_j\| - \|x - \tilde{c}_i\| > \tilde{\rho}_i + \bar{\rho}$ , implying  $\|x - \tilde{c}_j\|^2 > \tilde{\rho}_j^2 + \bar{\rho}^2$  and consequently,  $L_j(\|x - \tilde{c}_j\|^2 - \tilde{\rho}_j^2) = 0$ . Moreover, it holds that  $\|\tilde{c}_\ell - \tilde{c}_i\| \geq \rho_\ell + \rho_i + \bar{\rho}$  which implies that  $\|\tilde{c}_\ell - \tilde{c}_i\|^2 \geq \tilde{\rho}_i^2 + \tilde{\rho}_\ell^2 + \bar{\rho}^2$ , for all  $i \in \mathcal{K} \setminus \{k, \ell\}$ . Consequently,  $\|\tilde{c}_\ell - \tilde{c}_i\|^2 - \tilde{\rho}_i^2 > \bar{\rho}^2$ , which leads to  $L_i(\|\tilde{c}_\ell - \tilde{c}_i\|^2 - \tilde{\rho}_i^2) = 0$  for all  $i \in \tilde{\mathcal{K}} \setminus \{k, \ell\}$ . Finally, the function  $L_{k,\ell} : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$  is defined as  $L_{k,\ell}(e_\ell) := \sum_{i \in \tilde{\mathcal{K}} \setminus \{k, \ell\}} L_i(\|e_\ell + \tilde{c}_\ell - \tilde{c}_i\|^2 - \tilde{\rho}_i^2)$ . One concludes from above that  $L_{k,\ell}(0) = 0$ .