# Learning self-triggered controllers with Gaussian processes

Kazumune Hashimoto, Yuichi Yoshimura, Toshimitsu Ushio

*Abstract*—**This paper investigates the design of self-triggered controllers for networked control systems (NCSs), where the dynamics of the plant is *unknown* apriori. To deal with the unknown transition dynamics, we employ the Gaussian process (GP) regression in order to learn the dynamics of the plant. To design the self-triggered controller, we formulate an optimal control problem, such that the optimal control and communication policies can be jointly designed based on the GP model of the plant. Moreover, we provide an overall implementation algorithm that jointly learns the dynamics of the plant and the self-triggered controller based on a reinforcement learning framework. Finally, a numerical simulation illustrates the effectiveness of the proposed approach.**

*Index Terms*—**Event-triggered/self-triggered control, Optimal control, Gaussian process regression.**

## I. INTRODUCTION

In networked control systems (NCSs), sensors, actuators, and controllers reside in multiple areas linked by wired/wireless communication network. Due to the progress in communication technology and many practical advantages such as a low-cost maintenance and flexibility for modifications, NCSs have been developed in a wide variety of applications, including manufacturing plants, autonomous robots/vehicles, traffic networks, to name a few [1]. In recent years, event-triggered and self-triggered control have attracted much attention and are known to be useful strategies for the NCSs [2]. This is due to the fact that, it leads to the potential saving of resources that are present in NCSs, such as a limited battery capacity or a limited communication bandwidth, by transmitting sensor measurements over the communication network only when it is needed. So far, various event/self-triggered controllers have been proposed in the literature, see, e.g., [3] for survey papers. Early works consider designing event/self-triggered control based on input-to-state stability (ISS) or $\mathcal{L}_2$-gain performance [4]–[6]. More recently, event-triggered control has been formulated as the hybrid dynamical systems [7], [8]. In addition, some approaches to combine event/self-triggered control and optimal control have been also provided in recent years [9]–[20].

In the aforecited event-triggered and self-triggered control framework, it is generally assumed that the transition dynamics, which represents the underlying model of the plant, is *known* apriori. This implies that, when the event/self-triggered

controllers are applied to the real world (actual) control systems, the resulting performance is heavily dependent on how the system model is accurate with respect to the true dynamics. However, it may be the case in practice when an accurate model of the plant is hard to obtain based on the first principles from physics, due to the fact that the dynamics is complex and highly nonlinear. Examples include mechanical systems [21], autonomous vehicles [22], power consumption of multi-story buildings [23], periodic errors in astrophotography systems [24], to name a few.

Motivated by the above, in this paper we investigate the design of a novel self-triggered controller for NCSs, where the dynamics of the plant is assumed to be *unknown* apriori. To this end, we make use of the Gaussian process (GP) regression [25] in order to learn the dynamics of the plant. The use of GP offers many benefits, such as the ability to incorporate prior knowledge about the model (e.g., smoothness, periodicity) by selecting suitable kernel functions, as well as the ability to provide uncertainty of the model for prediction values. To design the self-triggered controller, we first formulate an infinite horizon optimal control problem, such that both the cost for the control performance and the communication are taken into account. Then, we derive the corresponding Bellman equation and provide an approach to solving the optimal control problem, such that both the optimal control and communication policies are designed based on the plant learned by the GP regression. In particular, we employ a value iteration algorithm, which derives the optimal policies by iteratively improving the estimate of the optimal cost function. Moreover, when solving the value iteration algorithm, we employ the so-called *moment matching* technique in order to approximate the multiple-ahead predictive distribution of states by the Gaussian distribution. As we will see later, this approximation together with the approximations of the optimal cost function based on the radial basis functions will allow us to derive the optimal policies in a tractable way. Finally, we provide an overall implementation algorithm that jointly learns the dynamics of the plant as well as the optimal control and the communication policies based on a reinforcement learning framework. As we will see later, this algorithm combines the *exploration/exploitation phase* that aims at collecting the training data to learn the dynamics of the plant in an $\varepsilon$-greedy fashion, and the *learning phase* that aims at updating the optimal control and communication policies based on the value iteration algorithm.

In summary, the main contributions of this paper is provided as follows:

1) We formulate an infinite horizon optimal control prob-

lem, such that both the control and communication policies can be designed based on the GP model of the plant.

2) We derive the Bellman equation corresponding to the optimal control problem and employ the value iteration algorithm to solve it. When solving this algorithm, we employ some approximation techniques, such as the moment matching, so that the (approximate) optimal policies can be derived.

3) We provide an overall reinforcement learning algorithm that jointly learns the GP model of the plant as well as the optimal control and communication policies.

*(Related works):* Our approach is related to several techniques that have been provided in the literature. Using the GP in control community has been attracted much attention in recent years [21]–[24], [26]–[29]. In particular, our approach is related to the ones based on optimal control framework, see, e.g., [22]–[24], [27]–[31]. For example, in [22], the authors have utilized the GP model to learn the dynamics of the plant, and they have formulated a chance-constrained model predictive control (MPC), in which the optimal control problem is solved for each time step based on the knowledge about the dynamics learned by the GP. In contrast to these previous methods, we provide an approach that jointly learns the dynamics of the plant and the *self-triggered controller*, aiming at reducing the number of communication time steps for NCSs. As previously mentioned and will be clearer in later sections, this is achieved by formulating a value iteration algorithm, such that the optimal pair of the control input and the inter-communication time steps can be determined for each state based on the GP dynamics of the plant.

With regard to the event/self-triggered control, some model-free/model-based approaches with unknown transition dynamics have been proposed in recent years, e.g., [32]–[42]. For example, in [36]–[41], an actor-critic based $Q$-learning algorithm was proposed to learn the intermittent feedback controller under the event-triggered policy, and closed-loop stability was rigorously shown. Our approach differs from those previous works, in the sense that; (i) we provide a model-based solution to the problem of learning self-triggered controllers based on the GP regression; (ii) while previous works aim at learning a controller based on a prescribed structure of the event-triggered condition (i.e., the event is triggered when the error between the actual state and the latest triggered state exceeds a certain threshold), our approach aims at learning both control and communication policies from scratch; (iii) while previous works deal with either linear or nonlinear input-affine systems, our approach is applicable to general nonlinear systems. Moreover, in [32], a deep reinforcement learning was proposed to learn the event-triggered controller, and, similarly to our approach, the communication policy was designed from scratch. One of the potential advantages over this previous work may be that, since our approach is a model-based approach that explicitly incorporates the knowledge about the dynamics, it may require much fewer number of iterative tasks to learn the desired policies. Such data-efficiency (see, e.g., [30]) is indeed illustrated in the simulation example in

Section VII, where we show that the desired policies can be learned within 10 episodes, while model-free approaches may typically require hundreds or thousands of iterative tasks to learn them.

***Notation.*** Throughout the paper, we make use of the following notations. Let $\mathbb{N}$, $\mathbb{N}_{\geq 0}$, $\mathbb{N}_{> 0}$, $\mathbb{N}_{a:b}$ be the set of integers, non-negative integers, positive integers, and the set of integers in the interval $[a, b]$, respectively. Let $\mathbb{R}$, $\mathbb{R}_{\geq 0}$, $\mathbb{R}_{> 0}$ be the set of reals, non-negative reals and positive reals, respectively. For a square matrix $\boldsymbol{Q}$, we use $\boldsymbol{Q} \succ 0$ to denote that $\boldsymbol{Q}$ is positive definite. Let $\mathrm{diag}(a_1, a_2, \ldots, a_N)$ be the diagonal matrix whose (diagonal) elements are given by $a_1, \ldots, a_N \in \mathbb{R}$. Moreover, let $\mathrm{Blkdiag}(A_1, A_2, \ldots, A_N)$ be the block diagonal matrix that consists of a set of matrices $A_1, \ldots, A_N$.

## II. PRELIMINARIES OF GAUSSIAN PROCESS REGRESSION

In this section, we provide some basic concepts and useful properties of the Gaussian process (GP) regression. Consider a nonlinear function $h : \mathbb{R}^n \to \mathbb{R}$ expressed as

$$y = h(\mathbf{x}) + \varepsilon, \tag{1}$$

where $\mathbf{x} \in \mathbb{R}^n$ is the input, $y \in \mathbb{R}$ is the output, and $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$ is the Gaussian distributed white noise. In the GP regression, we assume that the function $h$ follows the GP. That is, for every set of a finite (or possibly infinite) number of inputs $\mathbf{x}_i \in \mathbb{R}^n$, $i = 1, \ldots, N$, the joint probability of the corresponding set of outputs $\mathbf{y} = [y_1, y_2, \ldots, y_N]^\mathsf{T}$ follows the multivariate Gaussian distribution, i.e., $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{K})$, where $\boldsymbol{K} \in \mathbb{R}^{N \times N}$ is the covariance matrix and is characterized by $K_{ij} = \mathsf{k}(\mathbf{x}_i, \mathbf{x}_j)$, where $K_{ij}$ is the $(i, j)$-component of $\boldsymbol{K}$ and $\mathsf{k} : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}_{\geq 0}$ is the positive definite kernel function.

In this paper, we assume that the kernel function $\mathsf{k}$ is given by the squared exponential covariance function:

$$\mathsf{k}(\mathbf{x}_i, \mathbf{x}_j) = \alpha^2 \exp\left(-\frac{1}{2}(\mathbf{x}_i - \mathbf{x}_j)^\mathsf{T} \boldsymbol{\Lambda}^{-1} (\mathbf{x}_i - \mathbf{x}_j)\right), \tag{2}$$

where $\boldsymbol{\Lambda} = \mathrm{diag}\left(\lambda_1^2, \ldots, \lambda_N^2\right)$ and $\{\alpha, \lambda_1, \ldots \lambda_N\}$ are the hyper-parameters. For a given set of input-output training data $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$, the predictive distribution of the output for a new test input $\mathbf{x}$ follows the Gaussian distribution, i.e., $p(y|\mathbf{x}, \mathcal{D}) = \mathcal{N}(\mu(\mathbf{x}), \sigma(\mathbf{x}))$. Here the mean and the variance are given by

$$\mu(\mathbf{x}) = \mathbf{k}_*^\mathsf{T}(\mathbf{x})(\boldsymbol{K} + \sigma_\varepsilon^2 \boldsymbol{I})^{-1} \mathbf{y}, \tag{3}$$

$$\sigma(\mathbf{x}) = \mathsf{k}(\mathbf{x}, \mathbf{x}) - \mathbf{k}_*^\mathsf{T}(\mathbf{x})(\boldsymbol{K} + \sigma_\varepsilon^2 \boldsymbol{I})^{-1} \mathbf{k}_*(\mathbf{x}), \tag{4}$$

where $\mathbf{y} = [y_1, y_2, \ldots, y_N]^\mathsf{T}$ and

$$\mathbf{k}_*(\mathbf{x}) = [\mathsf{k}(\mathbf{x}, \mathbf{x}_1), \ldots, \mathsf{k}(\mathbf{x}, \mathbf{x}_N)]^\mathsf{T}. \tag{5}$$

Suitable selections of the hyper-parameters $\{\alpha, \lambda_1, \ldots \lambda_N\}$ are given by evidence maximization, see, e.g., [25]. For simplicity of presentation, we write $h \sim \mathcal{GP}$ if the function $h$ follows the GP.

## III. PROBLEM STATEMENT

In this section, we describe the dynamics of the plant,

Fig. 1. Networked Control System considered in this paper.

overview of the self-triggered controller, and define the cost function to be minimized.

### A. Dynamics

We consider a networked control system (NCS) illustrated in Fig. 1. As shown in the figure, the controller and the learning agent are connected to the plant over the communication network. Roughly speaking, the learning agent is responsible for learning the dynamics of the plant as well as the optimal control and communication policies. On the other hand, the controller is responsible for transmitting the control inputs to operate the plant based on the control and communication policies derived by the learning agent. This implementation will be formally given later in this paper. Throughout the paper, we assume that the communication network is ideal; it induces neither packet dropouts nor any network delays.

The dynamics of the plant is given by the following non-linear systems:

$$\mathbf{x}_{k+1} = \boldsymbol{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{u}_k \in U, \tag{6}$$

for all $k \in \mathbb{N}_{\geq 0}$, where $\mathbf{x}_k \in \mathbb{R}^{n_x}$ is the state, $\mathbf{u}_k \in \mathbb{R}^{n_u}$ is the control input, $U \subset \mathbb{R}^{n_u}$ is the set of control inputs, and $\boldsymbol{f} : \mathbb{R}^{n_x} \times \mathbb{R}^{n_u} \to \mathbb{R}^{n_x}$ is the transition dynamics that is assumed to be *unknown* apriori. While the transition dynamics is unknown, it is assumed here that the equilibrium point is known; without loss of generality, we assume that the origin has the equilibrium point, i.e., $\mathbf{0} = \boldsymbol{f}(\mathbf{0}, \mathbf{0})$. The control goal is to stabilize the system towards the origin.

Since $\boldsymbol{f}$ is unknown apriori, we consider that each component of the unknown function, i.e., $f_i$, $i \in \mathbb{N}_{1:n_x}$ ($\boldsymbol{f} = [f_1, f_2, \ldots, f_{n_x}]^\mathsf{T}$) is modeled by the GP regression. That is, $f_i$ is learned from the input-output training data $\mathcal{D}_i = \{\mathbf{X}, \mathbf{y}_i\}$, where

$$\mathbf{X} = \left[ \begin{bmatrix} \mathbf{x}_0^* \\ \mathbf{u}_0^* \end{bmatrix}, \begin{bmatrix} \mathbf{x}_1^* \\ \mathbf{u}_1^* \end{bmatrix}, \ldots, \begin{bmatrix} \mathbf{x}_{N-1}^* \\ \mathbf{u}_{N-1}^* \end{bmatrix} \right], \tag{7}$$

$$\mathbf{y}_i = [x_{i,1}^*, x_{i,2}^*, \ldots, x_{i,N}^*]^\mathsf{T}. \tag{8}$$

In (7) and (8), $N \in \mathbb{N}_{>0}$ denotes the number of training data points, $[\mathbf{x}_n^{*\mathsf{T}}, \mathbf{u}_n^{*\mathsf{T}}]$, $n \in \mathbb{N}_{1:N}$ are the training inputs following the dynamics (6) (i.e., $\mathbf{x}_{n+1}^* = \boldsymbol{f}(\mathbf{x}_n^*, \mathbf{u}_n^*)$, $n \in \mathbb{N}_{0:N-1}$), and $x_{i,n}^*$, $i \in \mathbb{N}_{1:n_x}$ is the $i$-th element of $\mathbf{x}_n^*$ as the set of training outputs. We denote by $\mathsf{k}_i(\cdot, \cdot)$, $\boldsymbol{K}_i$ and $\{\alpha_i, \lambda_{i,1}, \ldots, \lambda_{i,N}\}$ the kernel function, covariance matrix and the hyper-parameters for the GP model of $f_i$, respectively. Moreover, we denote by $\mu_i(\mathbf{x}, \mathbf{u})$, $\sigma_i(\mathbf{x}, \mathbf{u})$ the mean and the covariance for the GP model of $f_i$ with an arbitrary test input $\widetilde{\mathbf{x}} = [\mathbf{x}^\mathsf{T}, \mathbf{u}^\mathsf{T}]^\mathsf{T}$,

respectively, i.e.,

$$\mu_i(\mathbf{x}, \mathbf{u}) = \mathbf{k}_{*,i}^\mathsf{T}(\widetilde{\mathbf{x}})(\boldsymbol{K}_i + \sigma_\varepsilon^2 \boldsymbol{I})^{-1} \mathbf{y}_i, \tag{9}$$

$$\sigma_i(\mathbf{x}, \mathbf{u}) = \mathsf{k}_i(\widetilde{\mathbf{x}}, \widetilde{\mathbf{x}}) - \mathbf{k}_{*,i}^\mathsf{T}(\widetilde{\mathbf{x}})(\boldsymbol{K}_i + \sigma_\varepsilon^2 \boldsymbol{I})^{-1} \mathbf{k}_{*,i}(\widetilde{\mathbf{x}}), \tag{10}$$

where $\mathbf{k}_{*,i}(\mathbf{x}) = [\mathsf{k}_i(\mathbf{x}, \mathbf{x}_1), \ldots, \mathsf{k}_i(\mathbf{x}, \mathbf{x}_N)]^\mathsf{T}$. That is, letting $\widehat{f}_i$ be the GP model of $f_i$, we have

$$\widehat{f}_i(\mathbf{x}, \mathbf{u}) \sim \mathcal{N}\left(\mu_i(\mathbf{x}, \mathbf{u}), \sigma_i(\mathbf{x}, \mathbf{u})\right). \tag{11}$$

Then, the overall GP model for $\boldsymbol{f} = [f_1, f_2, \ldots, f_{n_x}]^\mathsf{T}$ is given by

$$\widehat{\boldsymbol{f}}(\mathbf{x}, \mathbf{u}) \sim \mathcal{N}\left(\boldsymbol{\mu}(\mathbf{x}, \mathbf{u}), \boldsymbol{\Sigma}(\mathbf{x}, \mathbf{u})\right), \tag{12}$$

where $\widehat{\boldsymbol{f}} = [\widehat{f}_1, \widehat{f}_2, \ldots, \widehat{f}_{n_x}]^\mathsf{T}$ and

$$\boldsymbol{\mu}(\mathbf{x}, \mathbf{u}) = [\mu_1(\mathbf{x}, \mathbf{u}), \ldots, \mu_{n_x}(\mathbf{x}, \mathbf{u})]^\mathsf{T}, \tag{13}$$

$$\boldsymbol{\Sigma}(\mathbf{x}, \mathbf{u}) = \mathrm{diag}\left(\sigma_1(\mathbf{x}, \mathbf{u}), \ldots, \sigma_{n_x}(\mathbf{x}, \mathbf{u})\right). \tag{14}$$

### B. Overview of the self-triggered controller

Let us now define the control and communication policies. First, let $k_i$, $i = 0, 1, 2, \ldots$ with $k_0 = 0$ and $k_{i+1} > k_i$, $\forall i \in \mathbb{N}_{\geq 0}$ be the communication time steps when the plant transmits the state $x_{k_i}$ to the controller. In addition, let $m_i \in \mathbb{N}_{>0}$, $i \in \mathbb{N}_{\geq 0}$ be the corresponding inter-communication time steps, i.e., $m_i = k_{i+1} - k_i$, $\forall i \in \mathbb{N}_{\geq 0}$. In this paper, we implement a *self-triggered controller* [2], aiming at reducing the number of communication time steps between the plant and the controller. That is, we aim at designing the (deterministic) policies $\pi = \{\pi_{\mathrm{inp}}, \pi_{\mathrm{com}}\}$, where

- $\pi_{\mathrm{inp}} : \mathbb{R}^{n_x} \to \mathbb{R}^{n_u}$ is the *control policy*, which is a mapping from the state to the corresponding control input;
- $\pi_{\mathrm{com}} : \mathbb{R}^{n_x} \to \mathbb{N}_{1:M}$ is the *communication policy*, which is the mapping from the state to the corresponding inter-communication time steps.

Here, $M \in \mathbb{N}_{>0}$ denotes the maximum inter-communication time step, which means that inter-communication time step does not exceed $M$. This parameter is a user-defined parameter and is chosen apriori in order to formulate the optimal control problem. The basic procedure of the self-triggered controller is summarized as follows: for each $k_i$, $i \in \mathbb{N}_{\geq 0}$,

[Step 1] the plant measures the state $\mathbf{x}_{k_i}$ and transmits $\mathbf{x}_{k_i}$ to the controller;

[Step 2] the controller computes the control input and the inter-communication time steps as $\mathbf{u}_{k_i} = \pi_{\mathrm{inp}}(\mathbf{x}_{k_i})$ and $m_i = \pi_{\mathrm{com}}(\mathbf{x}_{k_i})$;

[Step 3] the controller transmits $\{\mathbf{u}_{k_i}, m_i\}$ to the plant, and the plant applies $\mathbf{u}_{k_i}$ constantly until the next communication time, i.e., $\mathbf{u}_k = \mathbf{u}_{k_i}$, $\forall k \in \mathbb{N}_{k_i, k_{i+1}-1}$, where $k_{i+1} = k_i + m_i$;

### C. Cost function to be minimized

In this paper, we consider the following infinite-horizon cost function to be minimized:

$$J^\pi(\mathbf{x}_{k_i}) = \sum_{\ell=i+1}^{\infty} \mathbb{E}_{\mathbf{x}_{k_\ell}}^\pi \left[ C_1(\mathbf{x}_{k_\ell}) + \gamma C_2(m_\ell) \right], \tag{15}$$

where $\mathbb{E}_{\mathbf{x}}^{\pi}[\cdot]$ denotes the expectation with respect to $\mathbf{x}$, $C_1 : \mathbb{R}^{n_x} \to \mathbb{R}_{\geq 0}$ represents the stage cost for the state, $C_2 : \mathbb{N}_{1:M} \to \mathbb{R}_{\geq 0}$ represents the communication cost that aims to penalize the inter-communication time steps, and $\gamma > 0$ is the weight associated to the communication cost. We assume that the cost for the state is characterized by polynomials or exponential functions. For example, exponential type of the cost function is given by

$$C_1(\mathbf{x}_{k_\ell}) = 1 - \exp\left\{-\frac{1}{2}\mathbf{x}_{k_\ell}^\mathsf{T} \mathbf{Q}\mathbf{x}_{k_\ell}\right\}, \qquad (16)$$

where $\mathbf{Q} \succ 0$ is a given positive definite matrix. Moreover, polynomial cost functions include quadratic type:

$$C_1(\mathbf{x}_{k_\ell}) = \mathbf{x}_{k_\ell}^\mathsf{T} \mathbf{Q}\mathbf{x}_{k_\ell}. \qquad (17)$$

As will be clearer in later sections, the above characterizations will allow us to provide analytical computations of the integrals with respect to the Gaussian probability distribution.

The communication cost is characterized as follows:

$$C_2(m_\ell) = M - m_\ell. \qquad (18)$$

Recall that $M$ is the maximum inter-communication time steps, i.e., $m_\ell \leq M, \forall \ell \in \mathbb{N}$. Hence, the total cost function defined in (15) aims at taking the cost of the control performance and the communication into account, and the parameter $\gamma$ regulates the trade-off between them. As will be formalized in later sections, we design the optimal control and communication policies $\pi = \{\pi_{\text{inp}}, \pi_{\text{com}}\}$, such that (15) is minimized. Note that, since the function $\mathbf{f}$ is unknown apriori and is learned by the GP regression, we will make use of the GP model $\widehat{\mathbf{f}}$ (see (12)) in order to derive the optimal solution; for details, see Section V.

**Remark 1** (On the case of $\gamma = 0$). Note that, even for the case $\gamma = 0$, communication reduction can be potentially achieved by minimizing (15). This is due to the fact that the total cost in (15) is defined by summing the stage costs *only at the communication time steps*, i.e., the cost will be accumulated only when the communication is given. Hence, reducing the number of communication leads to the reduction of the total cost, and, therefore, minimizing (15) leads to the communication reduction even for the case $\gamma = 0$. This interpretation will be also illustrated in the simulation example, where the communication reduction will be indeed achieved for the case $\gamma = 0$ in contrast to the time-triggered strategy; for details, see Section VII. □

## IV. APPROXIMATING MULTIPLE-AHEAD PREDICTIONS UNDER CONSTANT CONTROL INPUTS

In this section, we describe a way of how to approximate multiple-ahead predictions of states under constant control inputs, provided that the GP model of the plant is obtained. Suppose that, for given GP model $\widehat{\mathbf{f}}$ in (12) and a pair $(\mathbf{x}_k, \mathbf{u}) \in \mathbb{R}^{n_x} \times U$, we aim at computing the predictive distribution of the states with the constant control input $\mathbf{u}$, i.e., $p(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}), p(\mathbf{x}_{k+2}|\mathbf{x}_k, \mathbf{u}), \ldots$, where $\mathbf{x}_{k+m}$, $m \in \mathbb{N}_{>0}$ represent the state from $\mathbf{x}_k$ by applying $\mathbf{u}$ constantly for $m$

time steps. In this paper, we employ a moment matching technique [30] in order to approximate the predictive distributions by the Gaussian distribution. Since the functions $f_i$, $i \in \mathbb{N}_{1:n_x}$ are modeled by the GP, the predictive distribution of the state for $k+1$ is given by $p(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}) = \mathcal{N}(\boldsymbol{\mu}_{k+1}, \boldsymbol{\Sigma}_{k+1})$, where $\boldsymbol{\mu}_{k+1} = \boldsymbol{\mu}(\mathbf{x}_k, \mathbf{u})$, $\boldsymbol{\Sigma}_{k+1} = \boldsymbol{\Sigma}(\mathbf{x}_k, \mathbf{u})$ with

$$\boldsymbol{\mu}(\mathbf{x}_k, \mathbf{u}) = [\mu_1(\mathbf{x}_k, \mathbf{u}), \ldots, \mu_{n_x}(\mathbf{x}_k, \mathbf{u})]^\mathsf{T} \qquad (19)$$

$$\boldsymbol{\Sigma}(\mathbf{x}_k, \mathbf{u}) = \text{diag}\left(\sigma_1(\mathbf{x}_k, \mathbf{u}), \ldots, \sigma_{n_x}(\mathbf{x}_k, \mathbf{u})\right). \qquad (20)$$

Here, $\mu_i(\cdot)$, $\sigma_i(\cdot)$ ($i \in \mathbb{N}_{1:n_x}$) are given by (9) and (10), respectively. Now, suppose that we would like to compute the distribution of the predictive state for general $k + m$, $m = 2, 3, \ldots$. To this end, suppose that the predictive distribution of $\mathbf{x}_{k+\ell}$, $\ell \in \mathbb{N}_{1:m-1}$ is approximated by the Gaussian, i.e., $p(\mathbf{x}_{k+\ell}|\mathbf{x}_k, \mathbf{u}) \approx \mathcal{N}(\boldsymbol{\mu}_{k+\ell}, \boldsymbol{\Sigma}_{k+\ell})$. Then, the predictive distribution for $k + \ell + 1$ can be derived as follows:

$$\begin{aligned} &p(\mathbf{x}_{k+\ell+1}|\mathbf{x}_k, \mathbf{u}) \\ &= \int p(\widetilde{\mathbf{x}}_{k+\ell}|\mathbf{x}_k, \mathbf{u})p(\mathbf{x}_{k+\ell+1}|\widetilde{\mathbf{x}}_{k+\ell}, \mathbf{x}_k, \mathbf{u})\mathrm{d}\widetilde{\mathbf{x}}_{k+\ell}, \\ &= \int p(\widetilde{\mathbf{x}}_{k+\ell}|\mathbf{x}_k, \mathbf{u})p(\mathbf{x}_{k+\ell+1}|\widetilde{\mathbf{x}}_{k+\ell})\mathrm{d}\widetilde{\mathbf{x}}_{k+\ell}, \end{aligned} \qquad (21)$$

where we let $\widetilde{\mathbf{x}}_{k+\ell} = [\mathbf{x}_{k+\ell}^\mathsf{T}, \mathbf{u}_{k+\ell}^\mathsf{T}]^\mathsf{T}$ and $\mathbf{u}_{k+\ell}$ denotes the control input applied at $k+\ell$. Since the analytical computation of the integral in (21) cannot be given, we compute the mean and the covariance of the right hand side of (21) and approximate $p(\mathbf{x}_{k+\ell+1}|\mathbf{x}_k, \mathbf{u})$ by the Gaussian distribution. The integral in (21) involves the joint distribution $p(\widetilde{\mathbf{x}}_{k+\ell}|\mathbf{x}_k, \mathbf{u})$, which is further computed as

$$\begin{aligned} p(\widetilde{\mathbf{x}}_{k+\ell}|\mathbf{x}_k, \mathbf{u}) &= p(\mathbf{x}_{k+\ell}, \mathbf{u}_{k+\ell}|\mathbf{x}_k, \mathbf{u}) \\ &= p(\mathbf{x}_{k+\ell}|\mathbf{x}_k, \mathbf{u})p(\mathbf{u}_{k+\ell}|\mathbf{u}, \mathbf{x}_{k+\ell}) \end{aligned}$$

Since $\mathbf{u}$ is applied constantly, it follows that $\mathbf{u}_{k+\ell} = \mathbf{u}$, i.e., $p(\mathbf{u}_{k+\ell}|\mathbf{u}, \mathbf{x}_{k+\ell}) = p(\mathbf{u}_{k+\ell}|\mathbf{u}) = \text{Dirac}(\mathbf{u}_{k+\ell} - \mathbf{u})$, where $\text{Dirac}(\cdot)$ denotes the Dirac delta function. Hence, (21) leads to

$$\begin{aligned} &p(\mathbf{x}_{k+\ell+1}|\mathbf{x}_k, \mathbf{u}) \\ &= \int p(\mathbf{x}_{k+\ell}|\mathbf{x}_k, \mathbf{u})p(\mathbf{u}_{k+\ell}|\mathbf{u})p(\mathbf{x}_{k+\ell+1}|\widetilde{\mathbf{x}}_{k+\ell})\mathrm{d}\widetilde{\mathbf{x}}_{k+\ell} \\ &= \int p(\mathbf{x}_{k+\ell}|\mathbf{x}_k, \mathbf{u})p(\mathbf{x}_{k+\ell+1}|\mathbf{x}_{k+\ell}, \mathbf{u})\mathrm{d}\mathbf{x}_{k+\ell}, \end{aligned} \qquad (22)$$

where $p(\mathbf{x}_{k+\ell}|\mathbf{x}_k, \mathbf{u}) \approx \mathcal{N}(\boldsymbol{\mu}_{k+\ell}, \boldsymbol{\Sigma}_{k+\ell})$. Moreover, using the GP model in (12), we have $p(\mathbf{x}_{k+\ell+1}|\mathbf{x}_{k+\ell}, \mathbf{u}) \approx p(\widehat{\mathbf{f}}(\mathbf{x}_{k+\ell}, \mathbf{u})|\mathbf{x}_{k+\ell}, \mathbf{u}) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_{k+\ell}, \mathbf{u}), \boldsymbol{\Sigma}(\mathbf{x}_{k+\ell}, \mathbf{u}))$, where

$$\boldsymbol{\mu}(\mathbf{x}_{k+\ell}, \mathbf{u}) = [\mu_1(\mathbf{x}_{k+\ell}, \mathbf{u}), \ldots, \mu_{n_x}(\mathbf{x}_{k+\ell}, \mathbf{u})]^\mathsf{T},$$

$$\boldsymbol{\Sigma}(\mathbf{x}_{k+\ell}, \mathbf{u}) = \text{diag}\left(\sigma_1(\mathbf{x}_{k+\ell}, \mathbf{u}), \ldots, \sigma_{n_x}(\mathbf{x}_{k+\ell}, \mathbf{u})\right).$$

In the above, $\mu_i(\cdot)$ and $\sigma_i(\cdot)$ ($i \in \mathbb{N}_{1:n_x}$) are computed according to (9) and (10), respectively. Based on the above, let us compute the mean and the covariance of the right hand side of (22). From (22), the mean of $p(\mathbf{x}_{k+\ell+1}|\mathbf{x}_k, \mathbf{u})$ is given

by

$$\boldsymbol{\mu}_{k+\ell+1} = \mathbb{E}_{\mathbf{x}_{k+\ell}}\left[\mathbb{E}_{\mathbf{x}_{k+\ell+1}}[\mathbf{x}_{k+\ell+1}|\mathbf{x}_{k+\ell},\mathbf{u}]\right]$$

$$= \int p(\mathbf{x}_{k+\ell}|\mathbf{x}_k,\mathbf{u})\boldsymbol{\mu}(\mathbf{x}_{k+\ell},\mathbf{u})\mathrm{d}\mathbf{x}_{k+\ell}$$

$$= \int \mathcal{N}(\boldsymbol{\mu}_{k+\ell},\boldsymbol{\Sigma}_{k+\ell})\boldsymbol{\mu}(\mathbf{x}_{k+\ell},\mathbf{u})\mathrm{d}\mathbf{x}_{k+\ell}. \qquad (23)$$

The integral in (23) can be computed analytically and is given by $\mu_{i,k+\ell+1} = \boldsymbol{\beta}_i^\mathsf{T}\boldsymbol{\eta}_i$, where $\mu_{i,k+\ell+1}$ denotes the $i$-th component of $\boldsymbol{\mu}_{k+\ell+1}$, $\boldsymbol{\beta}_i = (\boldsymbol{K}_i + \sigma_\varepsilon^2\boldsymbol{I})^{-1}\mathbf{y}_i$ and $\boldsymbol{\eta}_i = [\eta_{i,1},\eta_{i,2},\ldots,\eta_{i,N}]$ with

$$\eta_{i,n} = \alpha_i^2\left|(\boldsymbol{\Lambda}_i)^{-1}\widetilde{\boldsymbol{\Sigma}}_{k+\ell} + \boldsymbol{I}\right|^{-1/2}$$
$$\times \exp\left(-\frac{1}{2}(\widetilde{\boldsymbol{\mu}}_{k+\ell} - \widetilde{\mathbf{x}}_n^*)^\mathsf{T}(\boldsymbol{\Lambda}_i + \widetilde{\boldsymbol{\Sigma}}_{k+\ell})^{-1}(\widetilde{\boldsymbol{\mu}}_{k+\ell} - \widetilde{\mathbf{x}}_n^*)\right),$$

for all $i \in \mathbb{N}_{1:n_x}$, $n \in \mathbb{N}_{1:N}$. In the above, we let $\widetilde{\boldsymbol{\mu}}_{k+\ell} = [\boldsymbol{\mu}_{k+\ell}^\mathsf{T},\mathbf{u}^\mathsf{T}]^\mathsf{T}$, $\widetilde{\mathbf{x}}_n^* = [\mathbf{x}_n^{*\mathsf{T}},\mathbf{u}_n^{*\mathsf{T}}]^\mathsf{T}$ (recall that $\mathbf{x}_n^*$, $\mathbf{u}_n^*$ are the $n$-th training input defined in (7)), and $\widetilde{\boldsymbol{\Sigma}}_{k+\ell} = \text{Blkdiag}(\boldsymbol{\Sigma}_{k+\ell},\mathbf{0}_{n_u\times n_u})$ with $\mathbf{0}_{n_u\times n_u}$ being the $n_u \times n_u$ zero matrix. The covariance matrix $\boldsymbol{\Sigma}_{k+\ell+1}$ can be obtained by considering diagonal elements $\sigma_{i,k+\ell+1}$ and off-diagonal elements $\sigma_{ij,k+\ell+1}$, $i \neq j$ (see, e.g., [30]). The diagonal elements are given by

$$\sigma_{i,k+\ell+1} = \mathbb{E}_{\mathbf{x}_{k+\ell}}\left[\text{Var}_{\mathbf{x}_{k+\ell+1}}[x_{i,k+\ell+1}|\mathbf{x}_{k+\ell},\mathbf{u}]\right]$$
$$\quad + \text{Var}_{\mathbf{x}_{k+\ell}}\left[\mathbb{E}_{\mathbf{x}_{k+\ell+1}}[x_{i,k+\ell+1}|\mathbf{x}_{k+\ell},\mathbf{u}]\right]$$
$$= \mathbb{E}_{\mathbf{x}_{k+\ell}}\left[\text{Var}_{\mathbf{x}_{k+\ell+1}}[x_{i,k+\ell+1}|\mathbf{x}_{k+\ell},\mathbf{u}]\right]$$
$$\quad + \mathbb{E}_{\mathbf{x}_{k+\ell}}\left[\mathbb{E}_{\mathbf{x}_{k+\ell+1}}^2[x_{i,k+\ell+1}|\mathbf{x}_{k+\ell},\mathbf{u}]\right] - \mu_{i,k+\ell+1}^2$$
$$= \boldsymbol{\beta}_i^\mathsf{T}\boldsymbol{L}_i\boldsymbol{\beta}_i + \alpha_i^2 - \text{Tr}\left((\boldsymbol{K}_i + \sigma_\varepsilon^2\boldsymbol{I})^{-1}\boldsymbol{L}_i\right)$$
$$\quad + \sigma_\varepsilon^2 - \mu_{i,k+\ell+1}^2, \qquad (24)$$

where $\text{Var}_\mathbf{x}[\cdot]$ is the variance with respect to $\mathbf{x}$, $\boldsymbol{L}_i$ is the $N \times N$ matrix, whose $(p,q)$-component (denoted as $L_{i,pq}$) is given by

$$L_{i,pq} = |\boldsymbol{R}_i|^{-1/2}\,\mathsf{k}_i(\widetilde{\mathbf{x}}_p^*,\widetilde{\boldsymbol{\mu}}_{k+\ell})\mathsf{k}_i(\widetilde{\mathbf{x}}_q^*,\widetilde{\boldsymbol{\mu}}_{k+\ell})$$
$$\times \exp\left(2\boldsymbol{\Lambda}_i^{-2}(\widetilde{\mathbf{x}}_{pq}^*)^\mathsf{T}(\widetilde{\boldsymbol{\Sigma}}_{k+\ell}^{-1} + 2\boldsymbol{\Lambda}_i^{-1})^{-1}\widetilde{\mathbf{x}}_{pq}^*\right),$$

where $\widetilde{\mathbf{x}}_{pq}^* = \frac{1}{2}(\widetilde{\mathbf{x}}_p^* + \widetilde{\mathbf{x}}_q^*) - \widetilde{\boldsymbol{\mu}}_{k+\ell}$, $\boldsymbol{R}_i = 2\boldsymbol{\Lambda}_i^{-1}\widetilde{\boldsymbol{\Sigma}}_{k+\ell} + \boldsymbol{I}$. The off-diagonal elements are given by

$$\sigma_{ij,k+\ell+1} = \boldsymbol{\beta}_i^\mathsf{T}\boldsymbol{L}_{ij}\boldsymbol{\beta}_j - \mu_{i,k+\ell+1}\mu_{j,k+m+1}, \qquad (25)$$

where $\boldsymbol{L}_{ij}$ is the $N \times N$ matrix, whose $(p,q)$-component (denoted as $L_{ij,pq}$) is given by

$$L_{ij,pq} = |\boldsymbol{R}_{ij}|^{-1/2}\,\mathsf{k}_i(\widetilde{\mathbf{x}}_p^*,\widetilde{\boldsymbol{\mu}}_{k+\ell})\mathsf{k}_j(\widetilde{\mathbf{x}}_q^*,\widetilde{\boldsymbol{\mu}}_{k+\ell})$$
$$\times \exp\left(-\frac{1}{2}\left(\widetilde{\mathbf{x}}_{pq,ij}^*\right)^\mathsf{T}\left(\boldsymbol{\Lambda}_i^{-1} + \boldsymbol{\Lambda}_j^{-1} + \widetilde{\boldsymbol{\Sigma}}_{k+\ell}^{-1}\right)^{-1}\widetilde{\mathbf{x}}_{pq,ij}^*\right),$$

where $\boldsymbol{R}_{ij} = (\boldsymbol{\Lambda}_i^{-1} + \boldsymbol{\Lambda}_j^{-1})\widetilde{\boldsymbol{\Sigma}}_{k+\ell} + \boldsymbol{I}$ and

$$\widetilde{\mathbf{x}}_{pq,ij}^* = \boldsymbol{\Lambda}_i^{-1}(\widetilde{\mathbf{x}}_p^* - \widetilde{\boldsymbol{\mu}}_{k+\ell}) + \boldsymbol{\Lambda}_j^{-1}(\widetilde{\mathbf{x}}_q^* - \widetilde{\boldsymbol{\mu}}_{k+\ell}). \qquad (26)$$

Based on the above, we can approximate $p(\mathbf{x}_{k+\ell+1}|\mathbf{x}_k,\mathbf{u})$ by

the Gaussian distribution as

$$p(\mathbf{x}_{k+\ell+1}|\mathbf{x}_k,\mathbf{u}) \approx \mathcal{N}(\boldsymbol{\mu}_{k+\ell+1},\boldsymbol{\Sigma}_{k+\ell+1}). \qquad (27)$$

Hence, by recursively applying the above procedure for all $\ell = 1,\ldots,m-1$, we can approximate $p(\mathbf{x}_{k+m}|\mathbf{x}_k,\mathbf{u})$ by the Gaussian distribution.

## V. APPROXIMATE VALUE ITERATION

In this section, we provide an approach to deriving the optimal self-triggered controller that minimizes (15), provided the GP model of the plant (12) is obtained. Let $J^*(\mathbf{x}_{k_i}) = \min_\pi J^\pi(\mathbf{x}_{k_i})$. From (15), the corresponding optimal Bellman equation is given by

$$J^*(\mathbf{x}_{k_i})$$
$$= \min_{\mathbf{u}_{k_i},m_i}\left\{\mathbb{E}_{\mathbf{x}_{k_{i+1}}}\left[C(\mathbf{x}_{k_{i+1}},m_{i+1}) + J^*(\mathbf{x}_{k_{i+1}})\right]\right\}, \qquad (28)$$

where $C(\mathbf{x},m) = C_1(\mathbf{x}) + \gamma C_2(m)$. Since the state space $\mathbb{R}^{n_x}$ and the input space $U$ for the dynamics in (6) are both infinite, deriving an explicit solution to (28) is in general intractable. Thus, we derive an approximated solution to (28) by employing a finite number of representative points in the state space and the input space, which are denoted as $\mathbf{x}_{R,1},\mathbf{x}_{R,2},\ldots,\mathbf{x}_{R,N_X} \in \mathbb{R}^{n_x}$ and $\mathbf{u}_{R,1},\mathbf{u}_{R,2},\ldots,\mathbf{u}_{R,N_U} \in U$, respectively, with $N_X$ and $N_U$ being the number of representative points. These representative points may be selected as the grid points in a given bounded region of $\mathbb{R}^{n_x}$ as well as $\mathbb{R}^{n_u}$, so that they include the origin (as we aim at stabilizing the state towards the origin). For simplicity of presentation, we let $\boldsymbol{X}_R = \{\mathbf{x}_{R,0},\mathbf{x}_{R,1},\ldots,\mathbf{x}_{R,N_X}\}$, $\boldsymbol{U}_R = \{\mathbf{u}_{R,1},\mathbf{u}_{R,2},\ldots,\mathbf{u}_{R,N_U}\}$. The optimal cost function (denoted as $\widehat{J}^*$) and the optimal control policy (denoted as $\widehat{\pi}_{\text{inp}}^*$) are then approximated by the exponential Radial Basis Functions (RBFs):

$$\widehat{J}^*(\mathbf{x}) = \sum_{n=1}^{N_X} w_{J,n}\exp\left(-\frac{1}{2\sigma_J^2}\|\mathbf{x} - \mathbf{x}_{R,n}\|^2\right), \qquad (29)$$

$$\widehat{\pi}_{\text{inp}}^*(\mathbf{x}) = \sum_{n=1}^{N_X} w_{u,n}\exp\left(-\frac{1}{2\sigma_u^2}\|\mathbf{x} - \mathbf{x}_{R,n}\|^2\right), \qquad (30)$$

where $\{w_{J,n}\}_{n=1}^{N_X}$, $\{w_{u,n}\}_{n=1}^{N_X}$ are the weights and $\sigma_J$, $\sigma_u$ are the width of the RBFs for $\widehat{J}^*$ and $\widehat{\pi}_{\text{inp}}^*$, respectively, which are the hyper-parameters to be designed and will be updated during the algorithm. Moreover, the optimal communication policy is approximated by $\widehat{\pi}_{\text{com}}^*(\mathbf{x}) = [\![\widehat{\pi}_{\text{com}}'(\mathbf{x})]\!]$, where $[\![a]\!]$ denotes the closest positive integer to $a$ (i.e., $[\![a]\!] = \arg\min_j\{|a-j| : j \in \mathbb{N}_{>0}\}$ ) and

$$\widehat{\pi}_{\text{com}}'(\mathbf{x}) = \sum_{n=1}^{N_X} w_{c,n}\exp\left(-\frac{1}{2\sigma_c^2}\|\mathbf{x} - \mathbf{x}_{R,n}\|^2\right). \qquad (31)$$

Here, $\{w_{c,n}\}_{n=1}^{N_X}$ and $\sigma_c$ are the hyper-parameters to be updated.

The iterative procedure to solve (28) follows the so-called *value iteration* [43], which is summarized in Algorithm 1. As shown in the algorithm, for each $\mathbf{x} \in \boldsymbol{X}_R$, we compute

**Algorithm 1** Approximate value iteration to derive the self-triggered controller.

1: Initialize the hyper-parameters to represent $\widehat{\pi}_{\text{com}}^*$, $\widehat{\pi}_{\text{inp}}^*$ and $\widehat{J}^*$;
2: **for** Iteration $= 1 : N_{\text{ite}}$ **do**
3:   **for** all $\mathbf{x} \in \boldsymbol{X}_R$ **do**
4:     $D^*(\mathbf{x}) \leftarrow \infty$;
5:     **for** all $(\mathbf{u}, m) \in \boldsymbol{U}_R \times \mathbb{N}_{1:M}$ **do**
6:       Compute $D(\mathbf{x}, \mathbf{u}, m)$ as follows:

$$D(\mathbf{x}, \mathbf{u}, m)$$
$$\leftarrow \mathbb{E}_{\mathbf{x}_m}\left[C_1(\mathbf{x}_m) + \gamma C_2(m') + \widehat{J}^*(\mathbf{x}_m)\right];$$

7:       **if** $D(\mathbf{x}, \mathbf{u}, m) < D^*(\mathbf{x})$ **then**
8:         $D^*(\mathbf{x}) \leftarrow D(\mathbf{x}, \mathbf{u}, m)$;
9:         $\mathbf{u}^*(\mathbf{x}) \leftarrow \mathbf{u}$;
10:        $m^*(\mathbf{x}) \leftarrow m$;
11:       **end if**
12:     **end for**
13:   **end for**
14:   Update the hyper-parameters to represent $\widehat{\pi}_{\text{com}}'$, $\widehat{\pi}_{\text{inp}}^*$ and $\widehat{J}^*$ by using the new training data:

$$\{\mathbf{x}_{R,n}, D^*(\mathbf{x}_{R,n})\}_{n=1}^{N_X}, \{\mathbf{x}_{R,n}, \mathbf{u}^*(\mathbf{x}_{R,n})\}_{n=1}^{N_X},$$
$$\{\mathbf{x}_{R,n}, m^*(\mathbf{x}_{R,n})\}_{n=1}^{N_X}; \tag{32}$$

15: **end for**

---

$D(\mathbf{x}, \mathbf{u}, m)$ for all $\mathbf{u} \in \boldsymbol{U}_R$ and $m \in \mathbb{N}_{1:M}$, which are specifically defined as

$$D(\mathbf{x}, \mathbf{u}, m) = \mathbb{E}_{\mathbf{x}_m}\left[C_1(\mathbf{x}_m) + \gamma C_2(m') + \widehat{J}^*(\mathbf{x}_m)\right]$$

$$= \int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})C_1(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m \tag{33}$$

$$+ \gamma \int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})C_2(m')\mathrm{d}\mathbf{x}_m \tag{34}$$

$$+ \int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})\widehat{J}^*(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m \tag{35}$$

where $\mathbf{x}_m$ is the state that is reached from $\mathbf{x}$ by applying $\mathbf{u}$ constantly for $m$ time steps, $m'$ is the inter-communication time steps determined for the state $\mathbf{x}_m$, i.e., $m' = \widehat{\pi}_{\text{com}}^*(\mathbf{x}_m)$. As shown in (33)–(35), it is required to compute the distribution $p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})$, as well as the three expected values (integrals) with respect to this distribution. In what follows, we provide a detailed way of computing these three terms.

*(Computation of $p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})$):* The term $p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})$ is the predictive distribution of the state from $\mathbf{x}$ by applying $\mathbf{u}$ constantly for $m$ time steps, which can be indeed approximated by the moment matching technique as discussed in Section IV. That is, we can approximate the distribution as $p(\mathbf{x}_m|\mathbf{x}, \mathbf{u}) \approx \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$, where $\boldsymbol{\mu}_m$ and $\boldsymbol{\Sigma}_m$ denote the mean and the covariance of $p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})$ that are computed by following the technique described in Section IV.

*(Computation of (33)):* Using the Gaussian approximation of $p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})$, the first term (33) is given by

$$\int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})C_1(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m$$

$$\approx \int \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)C_1(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m. \tag{36}$$

Since we assume that $C_1$ is characterized by polynomials or exponential, we can analytically compute the integral in (36). For example, if $C_1$ is given by (16), the integral in (36) further leads to

$$\int \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)C_1(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m$$

$$= 1 - \int \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)\exp\left\{-\frac{1}{2}\mathbf{x}_m^{\mathsf{T}}\boldsymbol{Q}\mathbf{x}_m\right\}\mathrm{d}\mathbf{x}_m$$

$$= 1 - \delta(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

where $\delta(\cdot)$ is given by $\delta(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) = |\boldsymbol{I} + \boldsymbol{\Sigma}_m\boldsymbol{Q}|^{-\frac{1}{2}}$ $\exp\left(-\frac{1}{2}\boldsymbol{\mu}_m\boldsymbol{Q}(\boldsymbol{I} + \boldsymbol{\Sigma}_m\boldsymbol{Q})^{-1}\boldsymbol{\mu}_m\right)$.

*(Computation of (34)):* The second term (34) can be computed as

$$\int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})C_2(m')\mathrm{d}\mathbf{x}_m$$

$$= \int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})(M - \widehat{\pi}_{\text{com}}^*(\mathbf{x}_m))\mathrm{d}\mathbf{x}_m$$

$$= M - \int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})\widehat{\pi}_{\text{com}}^*(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m. \tag{37}$$

which requires to compute $\int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})\widehat{\pi}_{\text{com}}^*(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m$. Using (31), we approximate this term as follows:

$$\int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})\widehat{\pi}_{\text{com}}^*(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m$$

$$\approx \int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})\widehat{\pi}_{\text{com}}'(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m$$

$$= \sum_{n=1}^{N_X} w_{c,n}\delta_{c,n}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) \tag{38}$$

where $\delta_{c,s}(\cdot)$ is given by $\delta_{c,s}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) = |\boldsymbol{I} + \sigma_c^{-2}\boldsymbol{\Sigma}_m|^{-\frac{1}{2}}$ $\exp\left(-\frac{1}{2\sigma_c^2}(\boldsymbol{\mu}_m - \mathbf{x}_{R,n})^{\mathsf{T}}(\boldsymbol{I} + \boldsymbol{\Sigma}_m\sigma_c^{-2})^{-1}(\boldsymbol{\mu}_m - \mathbf{x}_{R,n})\right)$.

*(Computation of (35)):* The third integral (35) can be approximated in a similar manner to the computation of (34). From (29) and using $p(\mathbf{x}_m|\mathbf{x}, \mathbf{u}) \approx \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$, we have

$$\int p(\mathbf{x}_m|\mathbf{x}, \mathbf{u})\widehat{J}^*(\mathbf{x}_m)\mathrm{d}\mathbf{x}_m = \sum_{n=1}^{N_X} w_{J,n}\delta_{J,n}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m), \tag{39}$$

where $\delta_{J,n}(\cdot)$ is given by $\delta_{J,n}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) = |\boldsymbol{I} + \sigma_J^{-2}\boldsymbol{\Sigma}_m|$ $\exp\left(-\frac{1}{2\sigma_J^2}(\boldsymbol{\mu}_m - \mathbf{x}_{R,n})^{\mathsf{T}}(\boldsymbol{I} + \boldsymbol{\Sigma}_m\sigma_J^{-2})^{-1}(\boldsymbol{\mu}_m - \mathbf{x}_{R,n})\right)$.

As shown in the algorithm (line 3–line 10), for each $\mathbf{x} \in \boldsymbol{X}_R$ we pick the smallest value among $D(\mathbf{x}, \mathbf{u}, m)$, $\mathbf{u} \in \boldsymbol{U}_R$, $m \in \mathbb{N}_{1:M}$, as well as the corresponding pair of the control input and inter-communication time steps, which we denote by $D^*(\mathbf{x})$, $\mathbf{u}^*(\mathbf{x})$, and $m^*(\mathbf{x})$, respectively. Consequently, we obtain $\{D^*(\mathbf{x}_{R,n}), \mathbf{u}^*(\mathbf{x}_{R,n}), m^*(\mathbf{x}_{R,n})\}_{n=1}^{N_X}$, and these are used as the new training data to update the hyper-parameters of $\widehat{J}^*$, $\widehat{\pi}_{\text{inp}}^*$, $\widehat{\pi}_{\text{com}}'$ in (29), (30), (31). For example, $\widehat{J}^*$ is updated by using the training data $\{\mathbf{x}_{R,n}, D^*(\mathbf{x}_{R,n})\}_{n=1}^{N_X}$, where $\mathbf{x}_{R,n}$, $n \in \mathbb{N}_{1:N_X}$ are the training inputs and $D^*(\mathbf{x}_{R,n})$, $n \in \mathbb{N}_{1:N_X}$ are the training outputs.

**Remark 2** (On the selection of $M$)**.** As $M$ is selected larger, we may achieve a more communication reduction, since the

Fig. 2. Flowchart of the overall algorithm (Algorithm 2). As shown in the figure, the algorithm mainly consists of the *exploration/exploitation phase* that aims at executing the self-triggered controller or collecting the training data to learn the dynamics of the plant, and the *learning phase* that aims at updating the GP model $\widehat{f}$ as well as the optimal control and communication policies based on the value iteration algorithm (Algorithm 1).

controller may increase the possibility to select larger inter-communication time steps. However, the execution time to derive the optimal policies may increase as $M$ is selected larger, due to the fact that the number of evaluations to compute $D(\mathbf{x}, \mathbf{u}, m)$ (line 6 in Algorithm 1) increases. Hence, the parameter $M$ may be carefully chosen by considering the tradeoff between the communication reduction for the NCS and the computation load to derive the optimal policies according to Algorithm 1. □

## VI. IMPLEMENTATION

In this section, we provide an overall implementation algorithm that jointly learns the dynamics of the plant and the self-triggered controller based on a reinforcement learning framework.

The overall algorithm is shown in Fig. 2 as a flowchart and the details are shown in Algorithm 2. Since we assume that the learning agent has no knowledge about the dynamics of the plant, we set the communication policy as $\widehat{\pi}_{\text{com}}^*(\mathbf{x}) \leftarrow 1$, $\forall \mathbf{x} \in \mathbb{R}^{n_x}$ (i.e., communication is given at every time step), so that the learning agent can efficiently collect the training data and learn the dynamics of the plant at the initial phase (line 1 in Algorithm 2). As shown in Fig. 2 and Algorithm 2, the algorithm mainly consists of the following two steps; exploration/exploitation phase (line 11–line 28 in Algorithm 2), and learning phase (line 30, line 31 in Algorithm 2). During the exploration/exploitation phase, the controller implements the self-triggered controller in an $\varepsilon$-greedy fashion, as well as updates the training data. In the algorithm, $\text{Uniform}(0,1)$ (line 12) generates a random real number from the interval $[0,1]$ according to the uniform distribution. That is, with the probability $\varepsilon$, a random control input with the one step

---

**Algorithm 2** Overall reinforcement learning algorithm.

**Input:** $\mathbf{x}_{\text{init}}$ (initial state), $N_{\text{epi}}$ (number of episodes), $\varepsilon \in [0,1)$ (threshold for the greedy policy);
**Output:** $\widehat{\pi}_{\text{inp}}^*$, $\widehat{\pi}_{\text{com}}^*$ (approximated optimal control and communication policies);

1:  Initialize the hyper-parameters to represent $\widehat{\pi}_{\text{inp}}^*$, and set $\widehat{\pi}_{\text{com}}^*(\mathbf{x}) \leftarrow 1$, $\forall \mathbf{x} \in \mathbb{R}^{n_x}$;
2:  $\mathbf{X} \leftarrow \{\}$, $\forall i \in \mathbb{N}_{1:n_x}$;
3:  $\mathbf{y}_i \leftarrow \{\}$, $\forall i \in \mathbb{N}_{1:n_x}$;
4:  $\mathcal{D}_i \leftarrow \{\mathbf{X}, \mathbf{y}_i\}$, $\forall i \in \mathbb{N}_{1:n_x}$ (initialize the training data);
5:  **for** $n_{\text{epi}} = 1 : N_{\text{epi}}$ **do**
6:      $\ell \leftarrow 0$;
7:      $k_\ell \leftarrow 0$;
8:      $\mathbf{x}_{k_\ell} = \mathbf{x}_{\text{init}}$;
9:      The plant transmits $\mathbf{x}_{k_\ell}$ to the controller;
10:     [Exploration/Exploitation phase]
11:     **for** $\ell = 0 : N_{\max} - 1$ **do**
12:         Sample $r \sim \text{Uniform}[0,1]$;
13:         **if** $r < \varepsilon$ **then**
14:             $m_\ell \leftarrow 1$;
15:             Select $\mathbf{u}_{k_\ell}$ randomly from $U$;
16:         **else**
17:             $\mathbf{u}_{k_\ell} \leftarrow \widehat{\pi}_{\text{inp}}^*(\mathbf{x}_{k_\ell})$;
18:             $m_\ell \leftarrow \widehat{\pi}_{\text{com}}^*(\mathbf{x}_{k_\ell})$;
19:         **end if**
20:     **end for**
21:     $k_{\ell+1} \leftarrow k_\ell + m_\ell$;
22:     The controller transmits $\{\mathbf{u}_{k_\ell}, m_\ell\}$ to the plant;
23:     The plant applies $\mathbf{u}_{k_\ell}$ constantly for $m_\ell$ time steps and transmit $\mathbf{x}_{k_{\ell+1}}$ to the controller;
24:     **if** $m_\ell = 1$ **then**
25:         $\mathbf{X} \leftarrow \{\mathbf{X} \cup [\mathbf{x}_{k_\ell}^\mathsf{T}, \mathbf{u}_{k_\ell}^\mathsf{T}]^\mathsf{T}\}$;
26:         $\mathbf{y}_i \leftarrow \{\mathbf{y}_i \cup x_{k_{\ell+1},i}\}$, $\forall i \in \mathbb{N}_{1:n_x}$;
27:         $\mathcal{D}_i \leftarrow \{\mathcal{D}_i \cup \{\mathbf{X}, \mathbf{y}_i\}\}$;
28:     **end if**
29:     [Learning phase]
30:     The learning agent learns the GP model of the plant by using the new training data $\mathcal{D} = \{\mathcal{D}_i\}_{i=1}^{n_x}$;
31:     The learning agent executes Algorithm 1 to update the (approximated) optimal policies $\widehat{\pi}_{\text{inp}}^*$, $\widehat{\pi}_{\text{com}}^*$;
32: **end for**

---

inter-communication time step is sampled, and, otherwise, the computed optimal control and communication policies are chosen to be executed. Here, the one step inter-communication time step is chosen (with the probability $\varepsilon$) so that the learning agent is able to utilize the consecutive states (i.e., $\mathbf{x}_{k_\ell}$, $\mathbf{x}_{k_\ell+1}$ with $k_{\ell+1} = k_\ell + 1$) to update the GP model of $\boldsymbol{f}$. In the learning phase, the learning agent utilizes the new training data $\mathcal{D}$ to update the GP model of the plant, and compute the (approximated) optimal control and communication policies according to Algorithm 1.

Finally, some remarks on the proposed algorithm are in order as follows:

**Remark 3** (On achieving closed-loop stability)**.** Proving closed-loop stability by the proposed approach (Algorithm 2) is indeed challenging due to the following reasons. First, since we include the cost of communication in (15), if $\gamma$ (i.e., the weight for the communication cost) is selected too large, the penalty of the communication is too emphasized and the convergence to the origin may not be guaranteed. Indeed,

this issue will be pointed out in the simulation result, where it is shown that, as $\gamma$ is selected larger, the state does not converge to origin (for details, see Section VII). The closed-loop stability may be achieved as $\gamma \to 0$. However, since we employ the GP model of the plant when solving the optimal control problem, we first need to show that the GP model of the plant is accurate enough with respect to the true (actual) dynamics. Since there is no theoretical result on the error bound between the GP model $\widehat{f}$ and the true one $f$, how much training data should be collected to obtain the accurate model may be in general unknown. Hence, even though there exists a self-triggered controller that stabilizes the actual system to the origin, such stabilization is not guaranteed according to the policies derived according to Algorithm 2. $\qquad\square$

**Remark 4.** The lack of providing theoretical proof on closed-loop stability may be the drawback of our approach with respect to some previous works of event-triggered control with unknown transition dynamics (see, e.g., [36]–[41]). Nevertheless, our approach is advantageous over these previous works, in the sense that our approach is applicable to *general* nonlinear systems, while previous works focus on only input-affine or linear systems. For example, the prescribed event-triggered condition may be difficult to characterize for general nonlinear systems based on the procedure presented in [38], due to the fact that the Hamilton-Jacobi-Bellman (HJB) equation under the event-triggered strategy is no longer characterized by (13) in [38]. In this paper, the self-triggered controller for general nonlinear systems can be designed by learning the dynamics based on the GP regression and deriving both the control and communication policies from scratch by implementing Algorithm 1. $\qquad\square$

## VII. SIMULATION RESULTS

In this section, we illustrate the effectiveness of the proposed approach through a simulation example. The simulation was conducted on Matlab 2016a under Windows 10, Intel(R) Core(TM) i7 4.20 GHz, 32 GB RAM. As a simulation example, we consider a control problem of an inverted pendulum, whose dynamics is governed by

$$x_{1,k+1} = x_{1,k} + \Delta x_{2,k} \tag{40}$$
$$x_{2,k+1} = x_{2,k} + \Delta(\sin x_{1,k} - x_{2,k} + u_k), \tag{41}$$

where $x_{1,k}$ and $x_{2,k}$ with $\mathbf{x}_k = [x_{1,k}; x_{2,k}]$ are the states that represent the angular position and the velocity of the mass, $u_k \in \mathbb{R}$ is the control input, and $\Delta = 0.2$ denotes the sampling time interval. Letting $\boldsymbol{f}(\mathbf{x}_k, u_k) = [f_1(\mathbf{x}_k, u_k), f_2(\mathbf{x}_k, u_k)]^{\mathsf{T}}$ with $f_1(\mathbf{x}_k, u_k) = x_{1,k} + \Delta x_{2,k}$ and $f_2(\mathbf{x}_k, u_k) = x_{2,k} + \Delta(\sin x_{1,k} - x_{2,k} + u_k)$, we obtain the discrete-time system as $\mathbf{x}_{k+1} = \boldsymbol{f}(\mathbf{x}_k, u_k)$, $k \in \mathbb{N}$. Note that the function $\boldsymbol{f}(\cdot, \cdot)$ is assumed to be unknown apriori and is thus learned by the GP regression. It is assumed that $U = [-1.5, 1.5]$ and the initial state is given by $\mathbf{x}_{\text{init}} = [x_{1,0}; x_{2,0}] = [1.0; 0.2]$. The maximum inter-communication time step is $M = 10$, and the representative points for the state space to solve (28) is selected by the uniform grid points in the set $X = [-1.5, 1.5] \times [-1.5, 1.5]$ with the interval 0.3, i.e., $X_R = [X]_{0.3}$. The representative points for the input space is given by $U_R = [U]_{0.3}$. We use

the exponential type for the stage cost in (16) with $\boldsymbol{Q} = \boldsymbol{I}_2$, and we set $\gamma = 0$ for the cost function in (15).

Fig. 3(a) illustrates the trajectories of the states by applying the self-triggered controller obtained by Algorithm 2 with Episode $= 1$ (red dotted) and 10 (blue solid). The figures illustrate that, while the state diverges at the initial learning phase, it is indeed stabilized towards the origin as the number of episode increases. The computed inter-communication time steps corresponding to the simulation result in Fig. 3(a) (Episode $= 10$) are illustrated in Fig. 3(b), which shows that the communication is given aperiodically according to the derived self-triggered controller. Fig. 3(c) illustrates the state trajectories by applying Algorithm 2 after Episode $= 10$ with different selections of $M$ ($M = 1, 10$). Note that, $M = 1$ corresponds to the case when communication is given at every time step, i.e., the time-triggered controller is implemented. The figure shows that the convergence of states for the case $M = 1$ seems to be faster than for the case $M = 10$, which is due to the fact that control inputs are updated at every time step when the time-triggered controller is implemented. On the other hand, the total number of communication instants required for the time interval $k \in [0, 100)$ is 100 for the case $M = 1$ (as it is the time-triggered implementation), while it is 27 for the case $M = 10$. This implies that employing the self-triggered controller achieves a significant communication reduction in contrast to the time-triggered strategy. Hence, the result shows that there exists a tradeoff between the communication reduction for the NCS and the convergence speed of states towards the origin, and such tradeoff may be regulated by tuning the parameter $M$.

To indicate the robustness of the derived self-triggered controller, we also illustrate in Fig. 4 several trajectories from different initial states around $\mathbf{x}_{\text{init}}$. The figure illustrates that the states are indeed stabilized to the origin regardless of the deviation of the initial states, showing the robustness of the self-triggered controller.

Finally, to analyze the effect of $\gamma$, we illustrate in Fig. 5(a) and Fig. 5(b) the resulting state trajectories under different selections of $\gamma$ ($\gamma = 0.01, 0.02, 0.03$), and the corresponding inter-communication time steps, respectively. Here, Algorithm 2 has been implemented for each $\gamma$ with 10 episodes ($N_{\text{epi}} = 10$). From Fig. 5(b), it is shown that larger inter-communication time steps are more likely to be selected as $\gamma$ is selected larger. This is due to the fact that, by selecting larger $\gamma$, it will penalize more for the communication cost. Note that, for the case $\gamma = 0.03$, the resulting state trajectory converges farther from the origin than for the other cases (while it achieves larger inter-communication time steps), which may be due to the fact that achieving large inter-communication time steps is too emphasized. Hence, similarly to the above, there exists a tradeoff between the communication reduction for the NCS and the convergence of states towards the origin, and such tradeoff may be regulated by tuning the parameter $\gamma$.

## VIII. CONCLUSION AND FUTURE WORK

In this paper, we investigate the self-triggered controller for NCSs with the unknown transition dynamics. To this end,

we use the GP to learn the dynamics of the plant. We first formulate an optimal control problem, such that both the cost for the control performance and the communication cost can be taken into account. Then, we illustrate that the optimal control problem can be solved via a value iteration algorithm, in which the optimal pair of the control input and the inter-communication time steps can be determined based on the GP model of the plant. Then, we provide overall reinforcement learning algorithm that jointly learns the dynamics of the plant as well as the self-triggered controller implemented by the learning agent. Finally, a numerical simulation is given to illustrate the effectiveness of the proposed approach.

Future work involves analyzing some theoretical issues (e.g., stability of the closed loop system, convergence property of the value iteration algorithm, etc.) for the GP dynamics of the plant. Moreover, providing some experiments to test the applicability of our approach should be investigated for our future work of research.

## REFERENCES

[1] J. P. Hespanha, P. Naghshtabrizi, and Y. Xu, "A survey of recent results in networked control systems," *Proceedings of the IEEE*, vol. 95, no. 1, pp. 138–162, 2007.

[2] W. P. M. H. Heemels, K. H. Johansson, and P. Tabuada, "An introduction to event-triggered and self-triggered control," in *Proceedings of the 51st IEEE Conference on Decision and Control (IEEE CDC)*, 2012, pp. 3270–3285.

[3] C. Peng and F. Li, "A survey on recent advances in event-triggered communication and control," *Information Sciences*, vol. 457, pp. 113–125, 2018.

[4] X. Wang and M. D. Lemmon, "Self-triggered feedback control systems with finite $\mathcal{L}_2$ gain stability," *IEEE Transactions on Automatic Control*, vol. 54, no. 3, pp. 452–467, 2009.

[5] M. C. F. Donkers and W. P. M. H. Heemels, "Output-based event-triggered control with guaranteed $\mathcal{L}_\infty$ gain and decentralized event-triggering," *IEEE Transactions on Automatic Control*, vol. 57, no. 6, pp. 1362–1376, 2011.

[6] M. Mazo Jr., A. Anta, and P. Tabuada, "An iss self-triggered implementation of linear controllers," *Automatica*, vol. 46, no. 8, pp. 1310–1314, 2010.

[7] V. S. Dolk, D. P. Borgers, and W. P. M. H. Heemels, "Output-based and decentralized dynamic event-triggered control with guaranteed $\mathcal{L}_p$-gain performance and zeno-freeness," *IEEE Transactions on Automatic Control*, vol. 62, no. 1, pp. 34–49, 2016.

[8] W. P. M. H. Heemels, M. C. F. Donkers, and A. R. Teel, "Periodic event-triggered control for linear systems," *IEEE Transactions on Automatic Control*, vol. 58, no. 4, pp. 847–861, 2013.

[9] A. Eqtami, D. V. Dimarogonas, and K. J. Kyriakopoulos, "Event-triggered control for discrete time systems," in *Proceedings of American Control Conference (ACC)*, 2010, pp. 4719–4724.

[10] K. Hashimoto, S. Adachi, and D. V. Dimarogonas, "Self-triggered model predictive control for nonlinear input-affine dynamical systems via adaptive control samples selection," *IEEE Transactions on Automatic Control*, vol. 62, no. 1, pp. 177–189, 2017.

[11] ——, "Energy-aware networked control systems under temporal logic specifications," in *Proceedings of the 57th IEEE Conference on Decision and Control (IEEE CDC)*, 2018.

[12] ——, "Event-triggered intermittent sampling for nonlinear model predictive control," *Automatica*, vol. 81, pp. 148–155, 2017.

[13] K. Hashimoto and D. V. Dimarogonas, "Synthesizing communication plans for reachability and safety specifications," *IEEE Transactions on Automatic Control*, vol. 65, no. 2, pp. 561–576, 2020.

[14] ——, "Resource-aware networked control systems under temporal logic specifications," *Discrete Event Dynamic Systems*, vol. 29, pp. 473–499, 2019.

[15] K. G. Vamvoudakis, A. Mojoodi, and H. Ferraz, "Event-triggered optimal tracking control of nonlinear systems," *The International Journal of Robust and Nonlinear Control*, vol. 27, no. 4, pp. 598–619, 2017.

[16] A. Heydari, "Optimal triggering of networked control systems," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 7, pp. 3011–3021, 2018.

[17] Y. C. Sun and G. H. Yang, "Robust event-triggered model predictive control for cyber-physical systems under denial-of-service attacks," *The International Journal of Robust and Nonlinear Control*, vol. 29, no. 14, pp. 4797–4811, 2019.

[18] D. Tolic, R. Fierro, and S. Ferrari, "Optimal self-triggering for nonlinear systems via approximate dynamic programming," in *Proceedings of 2012 IEEE International Conference on Control Applications*, 2012, pp. 879–884.

[19] C. Liu, H. Li, Y. Shi, and D. Xu, "Co-design of event trigger and feedback policy in robust model predictive control," *IEEE Transactions on Automatic Control*, 2019(to appear).

[20] C. Liu, J. Gao, H. Li, and D. Xu, "Aperiodic robust model predictive control for constrained continuous-time nonlinear systems: An event-triggered approach," *IEEE Transactions on Cybernetics*, vol. 4, no. 5, pp. 1397–1405, 2018.

[21] T. Beckers, D. Kulic, and S. Hirche, "Stable gaussian process based tracking control of euler-lagrange systems," *Automatica*, vol. 103, pp. 390–397, 2019.

[22] M. N. Z. L. Hewing, A. Liniger, "Cautious nmpc with gaussian process dynamics for autonomous miniature race cars," in *Proceedings of 2018 European Control Conference (ECC 2018)*, 2018.

[23] A. Jain, T. X. Nghiem, M. Morari, and R. Mangharam, "Learning and control using gaussian processes: towards bridging machine learning and controls for physical systems," in *Proceedings of the 9th ACM/IEEE International Conference on Cyber-Physical Systems (ICCPS 2018)*, 2018.

[24] E. D. Klenske, M. N. Zeilinger, B. Scholkopf, and P. Hennig, "Gaussian process-based predictive control for periodic error correction," *IEEE Transactions on Control Systems Technology*, vol. 24, no. 1, pp. 390–397, 2019.

[25] C. F. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*, The MIT Press, 2006.

[26] J. Umlauft, L. Pohler, and S. Hirche, "An uncertainty-based control lyapunov approach for control-affine systems modeled by gaussian process," *IEEE Control Systems Letters*, vol. 2, no. 3, pp. 483–488, 2018.

[27] J. Kocijan, R. M. Smith, C. E., and A. Girard, "Gaussian process model predictive control," in *Proceedings of the 2004 American Control Conference*, 2004.

[28] J. Umlauft, T. Beckers, and S. Hirche, "Scenario-based optimal control for gaussian process state space models," in *Proceedings of 2018 European Control Conference (ECC 2018)*, 2018.

[29] E. Bradford, L. Imsland, D. Zhang, and E. A. R. Chanona, "Stochastic data-driven model predictive control using gaussian processes," in *arxiv*, available online at https://arxiv.org/pdf/1908.01786.pdf.

[30] M. P. Deisenroth, D. Fox, and C. E. Rasmussen, "Gaussian processes for data-efficient learning in robotics and control," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 2, pp. 408–423, 2013.

[31] M. P. Deisenroth, C. E. Rasmussen, and J. Peters, "Gaussian process dynamic programming," *Neurocomputing*, vol. 72, no. 7–9, pp. 1508–1524, 2009.

[32] D. Baumann, J.-J. Zhu, G. Martius, and S. Trimpe, "Deep reinforcement learning for event-triggered control," in *Proceedings of 57th IEEE Conference on Decision and Control (IEEE CDC)*, 2018, pp. 943–950.

[33] K. E. Årzen, "A simple event-based pid controller," in *Proceedings of 14th IFAC World Congress*, 1999.

[34] D. Baumann, F. Solowjow, K. H. Johansson, and S. Trimpe, "Event-triggered pulse control with model learning (if necessary)," in *Proceedings of 2019 American Control Conference (ACC 2019)*, 2019, pp. 792–797.

[35] J. Beuchert, F. Solowjow, J. Raisch, S. Trinpe, and T. Seel, "Hierarchical event-triggered learning for cyclically excited systems with application to wireless sensor networks," *IEEE Control Systems Letters*, vol. 4, no. 1, pp. 103–108, 2019.

[36] K. G. Vamvoudakis and H. Ferraz, "Model-free event-triggered control algorithm for continuous-time linear systems with optimal performance," *Automatica*, vol. 87, pp. 412–420, 2018.

[37] X. Zhong, Z. Ni, H. He, X. Xu, and D. Zhao, "Event-triggered reinforcement learning approach for unknown nonlinear continuous-time system," in *Proceedings of 2014 International Joint Conference on Neural Networks*, 2014.

[38] X. Yang and H. He, "Adaptive critic designs for event-triggered robust control of nonlinear systems with unknown dynamics," *IEEE Transactions on Cybernetics*, vol. 49, no. 6, pp. 2255–2267, 2019.

[39] Y. Yang, K. G. Vamvoudakis, H. Ferraz, and H. Modares, "Dynamic intermittent $Q$-learning-based model-free suboptimal co-design of $\mathcal{L}_2$-stabilization," *The International Journal of Robust and Nonlinear Control*, vol. 29, no. 9, pp. 2673–2694, 2019.

[40] ——, "Dynamic intermittent $Q$-learning for systems with reduced bandwidth," in *Proceedings of 2018 IEEE Conference on Decision and Control (IEEE CDC)*, 2018, pp. 924–931.

[41] Y. Yang, H. Modares, K. G. Vamvoudakis, Y. Yin, and D. C. Wunsch, "Dynamic intermittent feedback design for $H_\infty$ containment control on a directed graph," *IEEE Transactions on Cybernetics*, 2019.

[42] Y. Yang, K. G. Vamvoudakis, H. Modares, W. He, Y. Yin, and D. C. Wunsch, "Safe intermittent reinforcement learning for nonlinear systems," in *Proceedings of 2018 IEEE Conference on Decision and Control (IEEE CDC)*, 2019.

[43] D. P. Bertsekas, *Dynamic Programming and Optimal Control*, Athena Scientific, Belmont, MA, 1996.

(a) State trajectories by implementing Algorithm 2 with $M = 10$.



(b) Corresponding inter-communication time steps (Episode = 10).



(c) State trajectories by implementing Algorithm 2 with $M = 1, 10$.

Fig. 3. Simulation results by applying Algorithm 2. Fig. 3(a) illustrates the state trajectories by applying Algorithm 2 with $M = 10$ after Episode = 1 (red dotted) and 10 (blue solid). Fig. 3(b) illustrates the corresponding inter-communication time steps for the case Episode = 10. Moreover, Fig. 3(c) illustrates the state trajectories by applying Algorithm 2 after Episode = 10 with different selections of $M$ ($M = 1, 10$). Note that, $M = 1$ corresponds to the case when communication is given at every time step (i.e., the time-triggered controller is implemented).

Fig. 4. State trajectories from random initial states by applying the derived self-triggered controller.



(a) State trajectories with $\gamma = 0, 0.01, 0.03$.



(b) Inter-communication time steps with $\gamma = 0, 0.01, 0.03$.

Fig. 5. Simulation results with different selections of $\gamma$.