# Simulator-Driven Deceptive Control
# via Path Integral Approach

Apurva Patil[1,*], Mustafa O. Karabag[2,*], Takashi Tanaka[3], Ufuk Topcu[3]

*Abstract*— We consider a setting where a supervisor delegates an agent to perform a certain control task, while the agent is incentivized to deviate from the given policy to achieve its own goal. In this work, we synthesize the optimal deceptive policies for an agent who attempts to hide its deviations from the supervisor's policy. We study the deception problem in the continuous-state discrete-time stochastic dynamics setting and, using motivations from hypothesis testing theory, formulate a Kullback-Leibler control problem for the synthesis of deceptive policies. This problem can be solved using backward dynamic programming in principle, which suffers from the curse of dimensionality. However, under the assumption of deterministic state dynamics, we show that the optimal deceptive actions can be generated using path integral control. This allows the agent to numerically compute the deceptive actions via Monte Carlo simulations. Since Monte Carlo simulations can be efficiently parallelized, our approach allows the agent to generate deceptive control actions online. We show that the proposed simulation-driven control approach asymptotically converges to the optimal control distribution.

## I. Introduction

We consider a deception problem between a supervisor and an agent. The supervisor delegates an agent to perform a certain task and provides a reference policy to be followed in a stochastic environment. The agent, on the other hand, aims to achieve a different task and may deviate from the reference policy to accomplish its own task. The agent uses a deceptive policy to hide its deviations from the reference policy. In this work, we synthesize the optimal policies for such a deceptive agent.

We formulate the agent's deception problem using motivations from hypothesis testing theory. We assume that the supervisor aims to detect whether the agent deviated from the reference policy by observing the state-action paths of the agent. On the flip side, the agent's goal is to employ a deceptive policy that achieves the agent's task and minimizes the detection rate of the supervisor. We design the agent's deceptive policy that minimizes the Kullback-Leibler (KL) divergence from the reference policy while achieving the agent's task. The use of KL divergence is motivated by the log-likelihood ratio test, which is the most powerful detection test for any given significance level [1]. Minimizing the KL divergence is equivalent to minimizing the expected log-likelihood ratio between distributions of the paths generated by the agent's deceptive policy and the reference policy. We

also note that due to the Bratagnolle-Huber inequality [2], for any statistical test employed by the supervisor, the sum of false positive and negative rates is lower bounded by a decreasing function of KL divergence between the agent's policy and the reference policy. Consequently, minimizing the KL divergence is a proxy for minimizing the detection rate of the supervisor. We represent the agent's task with a cost function and formulate the agent's objective function as a weighted sum of the cost function and the KL divergence.

We assume that the agent's environment follows discrete-time continuous-state dynamics. When the dynamics are linear, the supervisor's (stochastic) policies are Gaussian, and the cost functions are quadratic, minimizing a weighted sum of the cost function, and the KL divergence leads to solving a linear quadratic regulator problem. However, we consider a broader setting with potentially non-linear state dynamics, non-quadratic cost functions, and non-Gaussian reference policies. In this case, the agent's optimal deceptive policy does not necessarily admit a closed-form solution. While the agent's problem can be solved using backward dynamic programming, this approach suffers from the curse of dimensionality.

We show that, under the assumption of deterministic state dynamics, the optimal deceptive actions can be generated using path integral control without explicitly synthesizing a policy. In detail, we propose a two-step randomized algorithm for simulator-driven control for deception. At each time step, the algorithm first creates forward Monte Carlo samples of system paths under the reference policy. Then, the algorithm uses a cost-proportional weighted sampling method to generate a control input at that time step. We show that the proposed approach asymptotically converges to the optimal action distribution. Since Monte Carlo simulations can be efficiently parallelized, our approach allows the agent to generate the optimal deceptive actions online.

The contributions of this paper are threefold: 1) The work studies a problem of deception under supervisory control for continuous-state discrete-time stochastic systems. Given a reference policy, we formalize the synthesis of an optimal deceptive policy as a KL control problem and solve it using backward dynamic programming. 2) For the deterministic state dynamics, we propose a path-integral-based solution methodology for simulator-driven control. We develop an algorithm based on Monte Carlo sampling to numerically compute the optimal deceptive actions online. Furthermore, we show that the proposed approach asymptotically converges to the optimal control distribution of the deceptive agent. 3) We present a numerical example to validate the

* Indicates equal contribution. [1]Walker Department of Mechanical Engineering, University of Texas at Austin, `apurvapatil@utexas.edu`. [2]Department of Electrical and Computer Engineering, University of Texas at Austin, `karabag@utexas.edu`. [3]Department of Aerospace Engineering and Engineering Mechanics, University of Texas at Austin, `ttanaka@utexas.edu, utopcu@utexas.edu` .

derived simulator-driven control synthesis framework.

### A. Related Work

Deception naturally occurs in settings where two parties with conflicting objectives coexist. The example domains for deception include robotics [3], [4], supervisory control settings [5], [6], warfare [7], and cyber systems [8].

We formulate a deception problem motivated by hypothesis testing. This problem has been studied for fully observable Markov decision processes [5], partially observable Markov decision processes [6], and hidden Markov models [9]. Different from [5], [6], [9] that study discrete-state systems and directly solve an optimization problem for the synthesis of deceptive policies, we consider a nonlinear continuous-state system and provide a sampling-based solution for the synthesis of deceptive policies. In the security framework, [10], [11] study the detectability of an attacker in a stochastic control setting. Similar to our formulation, [10], [11] provide a KL divergence-based optimization problem. While we consider an agent whose goal is to optimize a different cost function from the supervisor, [10], [11] consider an attacker whose goal is to maximize the state estimation error of a controller.

KL divergence objective is also used in reinforcement learning [12], [13] to improve the learning performance and in KL control frameworks [14], [15] for the efficient computation of optimal policies. In [15], Ito et al. studied the KL control problem for nonlinear continuous-state space systems and characterized the optimal policies. Different from [15], we provide a randomized control algorithm based on path integral approach that converges to the optimal policy as the number of samples increases. Path integral control is a sampling-based algorithm employed to solve nonlinear stochastic optimal control problems numerically [16], [17], [18]. It allows the policy designer to compute the optimal control inputs online using Monte Carlo samples of system paths. The Monte Carlo simulations can be massively parallelized on GPUs, and thus the path integral approach is less susceptible to the curse of dimensionality [19].

### B. Notation

Let $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ be a measurable space where $\mathcal{X} \subseteq \mathbb{R}^n$ is a Borel set and $\mathcal{B}(\mathcal{X})$ is a Borel $\sigma$-algebra. Suppose $(\Omega, \mathcal{F}, \mathcal{P})$ is a probability space. An $(\mathcal{F}, \mathcal{B}(\mathcal{X}))$-measurable random variable $X$ is a function $X : \Omega \to \mathcal{X}$ whose probability distribution $P_X$ is defined by

$$P_X(B) = \mathcal{P}(X \in B) \quad \forall B \in \mathcal{B}(\mathcal{X})$$

$P_{X_2|X_1}(\cdot|\cdot) : \mathcal{B}(\mathcal{X}_2) \times \mathcal{X}_1 \to [0, 1]$ represents a stochastic kernel on $\mathcal{X}_2$ given $\mathcal{X}_1$. For simplicity, we write $P_X(dx)$ and $P_{X_2|X_1}(dx_2|x_1)$ as $P(dx)$ and $P(dx_2|x_1)$. If $P_1$ and $P_2$ are probability distributions on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ then, the Kullback-Leibler (KL) divergence from $P_1$ to $P_2$ is defined as

$$D(P_2\|P_1) = \int_{\mathcal{X}} \log \frac{dP_2}{dP_1}(x) P_2(dx)$$

if the Radon-Nikodym derivative $\frac{dP_2}{dP_1}$ exists, and $D(P_2\|P_1) = +\infty$ otherwise. Throughout this paper, we use the natural logarithm. Let $\mathcal{T} = \{0, 1, ..., T\}$ be the set of discrete time indices. A set of variables $\{x_0, x_1, \ldots, x_T\}$ is denoted by $x_{0:T}$ and a Cartesian product of sets $\mathcal{X}_0 \times \mathcal{X}_1 \times \ldots \times \mathcal{X}_T$ is denoted by $\mathcal{X}_{0:T}$. $P_{X_{0:T}}(dx_{0:T})$ denotes the joint probability distribution of random variables $X_0, X_1, \ldots, X_T$ on $(\mathcal{X}_{0:T}, \mathcal{B}(\mathcal{X}_{0:T}))$.

## II. PROBLEM FORMULATION

We consider a setting in which a supervisor contracts an agent to perform a certain task. Suppose the agent operates in a stochastic environment and follows discrete-time continuous-state dynamics. Let the state transition law of the agent be denoted by $P(dx_{t+1}|x_t, u_t) : \mathcal{B}(\mathcal{X}_{t+1}) \times \mathcal{X}_t \times \mathcal{U}_t \to [0, 1]$, where the random variables $X_t \in \mathcal{X}_t$ and $U_t \in \mathcal{U}_t$ represent the state and the control input of the system at time $t \in \mathcal{T}$. $\mathcal{X}_t$ and $\mathcal{U}_t$ are assumed to be Euclidean spaces with appropriate dimensions. Suppose a supervisor provides a (possibly stochastic) reference policy $\{R_{U_t|X_t}(\cdot|x_t)\}_{t=0}^{T-1}$ to the agent and expects the agent to follow the policy to accomplish a certain task. Here, $R_{U_t|X_t} : \mathcal{B}(\mathcal{U}_t) \times \mathcal{X}_t \to [0, 1]$ is a stochastic kernel on $\mathcal{U}_t$ given $\mathcal{X}_t$. The agent, on the other hand, aims to achieve a different task by minimizing the following cost function, which we henceforth call as *path cost*:

$$C_{0:T}(x_{0:T}, u_{0:T-1}) := \sum_{t=0}^{T-1} C_t(x_t, u_t) + C_T(x_T) \quad (1)$$

where $C_t(\cdot, \cdot) : \mathcal{X}_t \times \mathcal{U}_t \to \mathbb{R}$ for $t \in \mathcal{T}$ and $C_T(\cdot) : \mathcal{X}_T \to \mathbb{R}$ represent the stage costs and the terminal cost, respectively. In order to minimize the path cost (1), the agent designs its policy (possibly stochastic) $\{Q_{U_t|X_t}(\cdot|x_t)\}_{t=0}^{T-1}$ that may deviate from the reference policy $\{R_{U_t|X_t}(\cdot|x_t)\}_{t=0}^{T-1}$. The agent also attempts to be stealthy to hide its deviations from the supervisor. While the agent executes its policy $Q$, suppose the supervisor observes its state-action paths $\{x_{0:T}, u_{0:T-1}\}$, and uses a likelihood ratio test to detect whether the agent deviates from the reference policy. According to the Neyman–Pearson lemma, the likelihood-ratio test is optimal among all simple hypothesis tests for a given significance level [1]. In other words, we consider the worst-case scenario for the agent to be detected by the supervisor. Suppose the initial state $X_0 = x_0$ of the agent is known. We denote the joint probability distribution of the state-action paths induced via the reference policy by

$$R_{X_{0:T} \times U_{0:T-1}}(dx_{0:T} \times du_{0:T-1})$$
$$= \prod_{t=0}^{T-1} P(dx_{t+1}|x_t, u_t) R(du_t|x_t), \quad (2)$$

and the joint distribution induced via the agent's policy by

$$Q_{X_{0:T} \times U_{0:T-1}}(dx_{0:T} \times du_{0:T-1})$$
$$= \prod_{t=0}^{T-1} P(dx_{t+1}|x_t, u_t) Q(du_t|x_t). \quad (3)$$

Given a path $\{x_{0:T}, u_{0:T-1}\}$ that is randomly sampled under the agent's policy, the supervisor computes the log-likelihood

ratio (LLR)

$$\pi(x_{0:T}, u_{0:T-1}) = \log \frac{dQ_{X_{0:T} \times U_{0:T-1}}}{dR_{X_{0:T} \times U_{0:T-1}}}(x_{0:T}, u_{0:T-1}). \quad (4)$$

The supervisor decides that the agent uses the reference policy $R$ if $\pi(x_{0:T}, u_{0:T-1}) \leq c$, and deviates from $R$ otherwise. Here $c$ is a constant chosen by the supervisor to obtain a specified significance level. An agent not wanting to be detected by the supervisor must minimize the LLR (4). However, since the agent's trajectories are stochastic, the agent cannot directly minimize the LLR. We consequently consider that the agent's goal is to minimize the expected LLR as follows:

$$\Pi = \mathbb{E}_Q\left[\log \frac{dQ_{X_{0:T} \times U_{0:T-1}}}{dR_{X_{0:T} \times U_{0:T-1}}}(x_{0:T}, u_{0:T-1})\right] \quad (5)$$

where $\mathbb{E}_Q[\cdot]$ represents the expectation with respect to the probability distribution $Q$ (3). Note that equation (5) also defines the Kullback-Leibler (KL) divergence $D(Q\|R)$ between the agent's distribution $Q$ and the reference distribution $R$. It can be shown that $D(Q\|R)$ can be written as the stage-additive KL divergence between $Q_{U_t|X_t}$ and $R_{U_t|X_t}$ as follows (see Appendix A):

$$D(Q\|R) = \mathbb{E}_Q\left[\sum_{t=0}^{T-1} D(Q_{U_t|X_t}(\cdot|X_t)\|R_{U_t|X_t}(\cdot|X_t))\right]. \quad (6)$$

Since the KL divergence $D(Q\|R)$ is equivalent to the expected LLR (5), in this work, the KL divergence is used as a proxy for the measure of the agent's deviations from the reference policy.

Minimizing the KL divergence is in fact equivalent to minimizing the detection rate of an attacker for an ergodic process as proved in [11]. While we do not consider an ergodic process, the use of KL divergence is still well-motivated by the Bretagnolle–Huber inequality [2]. Let $\mathcal{E}$ be an arbitrary set of events that the supervisor will identify the agent as a deceptive agent, i.e., the agent followed $Q$. According to the Bretagnolle–Huber inequality, we have

$$\Pr(\mathcal{E}|R) + \Pr(\neg\mathcal{E}|Q) \geq \frac{1}{2}\exp(-D(Q\|R)) \quad (7)$$

where $\Pr(\mathcal{E}|R)$ and $\Pr(\neg\mathcal{E}|Q)$ denote the supervisor's false positive and negative rates, respectively. The false positive rate is the probability that the supervisor will identify the well-intentioned agent as a deceptive agent, i.e., the agent's policy is $R$, but the supervisor thinks that the agent has followed $Q$. Similarly, the false negative rate is the probability that the supervisor will identify the deceptive agent as a well-intentioned agent. The Bretagnolle–Huber inequality (7) states that the sum of the supervisor's false positive and negative rates is lower bounded by a decreasing function of the KL divergence between the distributions $Q$ and $R$. Therefore, an agent wanting to increase the supervisor's false classification rate should minimize the KL divergence from $R$ to $Q$.

The goal of the agent is to design a deceptive policy $Q$ that minimizes the expected path cost $\mathbb{E}_Q[C_{0:T}(X_{0:T}, U_{0:T-1})]$

(1) while limiting the KL divergence $D(Q\|R)$ (6). Using (1) and (6), we propose the following KL control problem for the synthesis of optimal deceptive policies for the agent:

**Problem 1** (Synthesis of optimal deceptive policy)**.**

$$\min_{\{Q_{U_t|X_t}\}_{t=0}^{T-1}} \mathbb{E}_Q \sum_{t=0}^{T-1}\left\{C_t(X_t, U_t) \quad (8)\right.$$

$$\left. + \lambda D(Q_{U_t|X_t}(\cdot|X_t)\|R_{U_t|X_t}(\cdot|X_t))\right\} + \mathbb{E}_Q C_T(X_T)$$

where $\lambda$ is a positive weighting factor that balances the trade-off between the KL divergence and the path cost.

We explain the above KL control problem via the following example:

**Example 1.** *Consider a drone that is contracted by a supervisor to perform a surveillance task over an area. The supervisor prefers the drone to operate at high speeds (policy $R$) to improve the efficiency of the surveillance. The operator of the drone, the agent, on the other hand, prefers the drone to operate in a battery-saving, safe mode (policy $Q$) to improve the longevity of the drone. The agent does not want to get detected by the supervisor and fired. Hence, the goal of the agent is to operate in a way that would balance the energy consumption ($\mathbb{E}_Q[C_{0:T}(X_{0:T}, U_{0:T-1})]$) and the deviations from the behavior desired by the supervisor ($D(Q\|R)$).*

## III. Synthesis of Optimal Deceptive Policies

In this section, we solve Problem 1 using backward dynamic programming and propose a policy synthesis algorithm based on path integral control.

### A. Backward Dynamic Programming

Notice that the cost function of Problem 1 possesses the time-additive Bellman structure and, therefore, can be solved by utilizing the principle of dynamic programming [20]. Define for each $t \in \mathcal{T}$ and $x_t \in \mathcal{X}_t$, the value function:

$$J_t(x_t) := \inf_{\{Q_{U_k|X_k}\}_{k=t}^{T-1}} \mathbb{E}_Q \sum_{k=t}^{T-1}\left\{C_k(X_k, U_k) \quad (9)\right.$$

$$\left. + \lambda D(Q_{U_k|X_k}(\cdot|X_k)\|R_{U_k|X_k}(\cdot|X_k))\right\} + \mathbb{E}_Q C_T(X_T).$$

Notice that in (9), we used "inf" instead of "min" since we do not know if the infimum is attained. In the following theorem, we show that the infimum is indeed attained, and therefore, "inf" can be replaced by "min".

**Theorem 1.** *The value function $J_t(x_t)$ satisfies the following backward Bellman recursion with the terminal condition $J_T(x_T) = C_T(x_T)$:*

$$J_t(x_t) = -\lambda \log\left\{\int_{\mathcal{U}_t} \exp\left(-\frac{C_t(x_t, u_t)}{\lambda}\right) \quad (10)\right.$$

$$\left.\times \exp\left(-\frac{1}{\lambda}\int_{\mathcal{X}_{t+1}} J_{t+1}(x_{t+1})P(dx_{t+1}|x_t, u_t)\right)R(du_t|x_t)\right\}$$

and the minimizer of Problem 1 is given by

$$Q^*_{U_t|X_t}(B_{U_t}|x_t) = \frac{\int_{B_{U_t}} \exp(-\rho_t(x_t, u_t)/\lambda) R(du_t|x_t)}{\int_{\mathcal{U}_t} \exp(-\rho_t(x_t, u_t)/\lambda) R(du_t|x_t)} \quad (11)$$

where

$$\rho_t(x_t, u_t) := C_t(x_t, u_t) + \int_{\mathcal{X}_{t+1}} J_{t+1}(x_{t+1}) P(dx_{t+1}|x_t, u_t) \quad (12)$$

and $B_{U_t}$ is a Borel set belonging to the $\sigma-$algebra $\mathcal{B}(\mathcal{U}_t)$.

*Proof.* See Appendix C. ∎

Theorem 1 provides a recursive method to compute the value functions $J_t(x_t)$ and optimal control distributions $Q^*_{U_t|X_t}$ backward in time. As one can see, to perform the backward recursions (10) and (11), the function $J_t(x_t)$ must be evaluated everywhere in the continuous domain $\mathcal{X}_t$. Therefore, in practice, an exact implementation of backward dynamic programming is computationally costly (unless the problem has a special structure, for example, linear state dynamics and quadratic costs). The computational cost grows quickly with the dimension of the state space of the system, which is referred to as the *curse of dimensionality*. In the next section (Section III-B), we show that under the assumption of the deterministic state transition law, the backward Bellman recursions can be linearized. This allows us to design a simulator-driven algorithm to compute optimal deceptive actions.

### B. Simulator-Driven Control via Path Integral Approach

In this section, we focus on a special case in which the agent's dynamics are deterministic and propose a simulator-driven algorithm to compute the optimal deceptive actions via path integral control.

**Assumption 1.** *The state transition law is governed by a deterministic mapping $F_t : \mathcal{X}_t \times \mathcal{U}_t \to \mathcal{X}_{t+1}$ as*

$$x_{t+1} = F_t(x_t, u_t); \quad (13)$$

*that is, $P(dx_{t+1}|x_t, u_t) = \delta_{F_t(x_t, u_t)}(dx_{t+1})$, where $\delta$ denotes the Dirac measure.*

**Remark 1.** *Note that under Assumption 1, the agent can deviate from the reference policy $R$ only if it is stochastic. Otherwise, under any deviations from the reference policy, with a positive probability, the supervisor will be sure that the agent did not follow the reference policy. Therefore, in what follows, we consider the reference policy to be stochastic. The stochasticity of the reference policy could be to account for the unmodeled elements of the dynamics, to provide robustness, or to encourage exploration.*

**Remark 2.** *Consider a special setting in which the state dynamics $F_t(x_t, u_t)$ is linear in $x_t$ and $u_t$, the reference policy distribution $R_{U_t|X_t}(\cdot|x_t)$ is Gaussian, and the cost functions $C_t(\cdot, \cdot)$ and $C_T(\cdot)$ are quadratic in $x_t$ and $u_t$. In such a setting, it can be shown that the optimal deceptive policy $Q^*_{U_t|X_t}(\cdot|x_t)$ is also Gaussian and can be analytically computed by backward Riccati recursions similar to the*

standard Linear-Quadratic-Regular (LQR) problems. In this work, we consider a broader setting with possibly non-Gaussian reference distribution, non-linear state dynamics, and non-quadratic cost functions. In this case, the optimal deceptive policy might not be efficiently computed by solving backward recursions.

Now, we propose a path-integral-based solution approach for simulator-driven policy synthesis. Under assumption 1, we can rewrite (10) as

$$J_t(x_t) = -\lambda \log \left\{ \int_{\mathcal{U}_t} \exp \left( -\frac{C_t(x_t, u_t)}{\lambda} \right) \quad (14) \right.$$
$$\left. \times \exp \left( -\frac{J_{t+1}(F_t(x_t, u_t))}{\lambda} \right) R(du_t|x_t) \right\}.$$

We introduce the exponentiated value function as $Z_t(x_t) := \exp\left(-\frac{1}{\lambda} J_t(x_t)\right)$. Using $Z_t(x_t)$, the Bellman recursion (14) can be linearized, and we get the following linear relationship between $Z_t$ and $Z_{t+1}$:

$$Z_t(x_t) = \int_{\mathcal{U}_t} \exp \left( -\frac{C_t(x_t, u_t)}{\lambda} \right) Z_{t+1}(F_t(x_t, u_t)) R(du_t|x_t)$$
$$= \int_{\mathcal{U}_t} \int_{\mathcal{X}_{t+1}} \exp \left( -\frac{C_t(x_t, u_t)}{\lambda} \right) Z_{t+1}(x_{t+1}) \quad (15)$$
$$\times P(dx_{t+1}|x_t, u_t) R(du_t|x_t).$$

Note that in (15), $P(dx_{t+1}|x_t, u_t) = \delta_{F_t(x_t, u_t)}(dx_{t+1})$ by Assumption 1. Equation (15) is a linear backward recursion in $Z_t$. The linear solvability of the KL control problem is well-known in the literature (e.g., [14]). We remark that linearizability critically relies on Assumption 1[1]. Now, by recursive substitution, (15) can also be written as

$$Z_t(x_t) = \int_{\mathcal{U}_t} \int_{\mathcal{X}_{t+1}} \cdots \int_{\mathcal{U}_{T-1}} \int_{\mathcal{X}_T} \exp \left( -\frac{C_t(x_t, u_t)}{\lambda} \right)$$
$$\times \cdots \times \exp \left( -\frac{C_T(x_T)}{\lambda} \right) R(dx_{t+1:T} \times du_{t:T-1}|x_t).$$

Thus, by introducing the path cost function

$$C_{t:T}(x_{t:T}, u_{t:T-1}) := \sum_{k=t}^{T-1} C_k(x_k, u_k) + C_T(x_T),$$

we obtain

$$Z_t(x_t) = \mathbb{E}_R \exp \left( -\frac{1}{\lambda} C_{t:T}(X_{t:T}, U_{t:T-1}) \right) \quad (16)$$

where the expectation $\mathbb{E}_R(\cdot)$ is with respect to the probability measure $R$ (2). Equation (16) expresses the exponentiated value function $Z_t(x_t)$ as the expected path cost under the reference distribution. This suggests a path-integral-based approach to numerically compute $Z_t(x_t)$. Suppose we generate a collection of $N$ independent samples of system paths $\{x_{t:T}(i), u_{t:T-1}(i)\}_{i=1}^N$ starting from $x_t$ under the reference

---

[1]We remark that in prior works where the path integral method is used to solve stochastic control problems, a certain assumption (e.g. Eq. (9) in [16]) is made to reinterpret the original problem as a problem of designing the optimal randomized policy for a deterministic transition system.

distribution $R$. Since the reference distribution is known, a collection of such sample paths can be easily generated using a Monte Carlo simulation. If $C_{t:T}(x_{t:T}(i), u_{t:T-1}(i))$ represents the path cost of the sample path $i$, then by the strong law of large numbers [21] as $N \to \infty$, we get

$$\frac{1}{N}\sum_{i=1}^{N} \exp\left(-\frac{1}{\lambda} C_{t:T}(x_{t:T}(i), u_{t:T-1}(i))\right) \overset{a.s.}{\to} Z_t(x_t). \quad (17)$$

Similarly, a collection of sample paths $\{x_{t:T}(i), u_{t:T-1}(i)\}_{i=1}^{N}$ starting from $x_t$ under the reference distribution $R$ can be used to sample $u_t$ from the optimal distribution $Q^*_{U_t|X_t}(\cdot|x_t)$. Notice that the optimal deceptive policy (11) can be expressed in terms of $\{Z_t\}_{t=0}^{T}$ as

$$
\begin{aligned}
Q^*_{U_t|X_t}(B_{U_t}|x_t) =& \frac{1}{Z_t(x_t)} \int_{B_{U_t}} \exp\left(-\frac{C_t(x_t, u_t)}{\lambda}\right) \\
& \times \exp\left(-\frac{J_{t+1}\left(F_t(x_t, u_t)\right)}{\lambda}\right) R(du_t|x_t) \\
=& \frac{1}{Z_t(x_t)} \int_{B_{U_t}} \exp\left(-\frac{C_t(x_t, u_t)}{\lambda}\right) \\
& \times Z_{t+1}\left(F_t(x_t, u_t)\right) R(du_t|x_t) \quad (18a) \\
=& \frac{1}{Z_t(x_t)} \int_{B_{U_t}} \int_{\mathcal{X}_{t+1}} \exp\left(-\frac{C_t(x_t, u_t)}{\lambda}\right) Z_{t+1}(x_{t+1}) \\
& \times P(dx_{t+1}|x_t, u_t) R(du_t|x_t) \quad (18b) \\
=& \frac{1}{Z_t(x_t)} \int_{\{\mathcal{X}_{t+1:T}, \mathcal{U}_{t:T-1}|u_t \in B_{U_t}\}} \exp\left(-\frac{C_{t:T}(x_{t:T}, u_{t:T-1})}{\lambda}\right) \\
& \times R(dx_{t+1:T} \times du_{t:T-1}|x_t). \quad (18c)
\end{aligned}
$$

The step (18a) follows from the definition of $Z_t$. In (18b), we used our assumption $P(dx_{t+1}|x_t, u_t) = \delta_{F_t(x_t, u_t)}(dx_{t+1})$. Finally, (18c) is obtained by the recursive substitution of (18b), and $\{\mathcal{X}_{t+1:T}, \mathcal{U}_{t:T-1}|u_t \in B_{U_t}\}$ represents a collection paths such that $u_t \in B_{U_t}$.

We use the above representation of $Q^*_{U_t|X_t}$ to sample an action $u_t$ from it. Let $r_t(i)$ be the exponentiated path cost of the sample path $i$:

$$r_t(i) := \exp\left(-\frac{1}{\lambda} C_{t:T}(x_{t:T}(i), u_{t:T-1}(i))\right) \quad (19)$$

and $r_t := \sum_{i=1}^{N} r_t(i)$. For each $t \in \mathcal{T}$, we introduce a piecewise constant, monotonically non-decreasing function $F_t : [0, N] \to [0, r_t]$ by

$$F_t(x) = \sum_{i=1}^{\lfloor x \rfloor} r_t(i).$$

where $\lfloor x \rfloor$ denotes $\text{floor}(x)$, i.e., the greatest integer less than or equal to $x$. The function $F_t(x)$ is represented in Figure 1.

Notice that the inverse $F_t^{-1}$ of $F_t$ defines a mapping $F_t^{-1} : [0, r_t] \to \{1, 2, ..., N\}$. To generate a sample $u_t$ approximately from the optimal distribution $Q^*_{U_t|X_t}$, we propose Algorithm 1. We first, generate a random variable $d$ according to $d \sim \text{unif}[0, r_t]$. Then, we select a sample ID by $j_t \leftarrow F_t^{-1}(d)$. Finally, the control input adopted in the $j_t$-th sample path at time step $t$ is selected as $u_t$, i.e.,
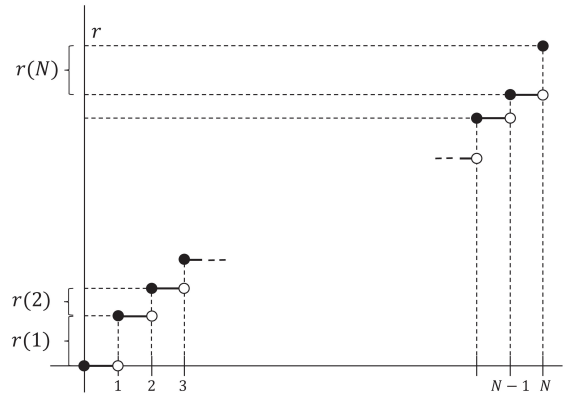


Fig. 1. Function $F_t(x)$.

---

**Algorithm 1:** Sampling $u_t$ approximately from $Q^*_{U_t|X_t}(\cdot|x_t)$ by Monte Carlo simulations

---
**Data:** Initial state $x_0$
1 **for** $t \in \mathcal{T}$ **do**
2      Sample $N$ paths $\{x_{t:T}(i), u_{t:T-1}(i)\}_{i=1}^{N}$ starting from $x_t$ under the reference distribution $R$.
3      For each sample path $i$, compute the exponentiated path cost $r_t(i)$ by (19).
4      Compute $r_t := \sum_{i=1}^{N} r_t(i)$.
5      Generate a random variable $d$ according to $d \sim \text{unif}[0, r_t]$.
6      Select a sample ID by $j_t \leftarrow F_t^{-1}(d)$.
7      Select a control input as $u_t \leftarrow u_t(j_t)$.

---

$u_t \leftarrow u_t(j_t)$. Theorem 2 proves that as the number of Monte Carlo samples tends to infinity, Algorithm 1 samples $u_t$ from the optimal distribution $Q^*_{U_t|X_t}(\cdot|x_t)$.

**Theorem 2.** *Let $B_{U_t} \in \mathcal{B}(\mathcal{U}_t)$ be a Borel set. Suppose for a given collection of sample paths $\{x_{t:T}(i), u_{t:T-1}(i)\}_{i=1}^{N}$, $u_t$ is computed by Algorithm 1 and the probability of $u_t \in B_{U_t}$ is denoted by $\Pr\{u_t \in B_{U_t}|\{x_{t:T}(i), u_{t:T-1}(i)\}_{i=1}^{N}\}$. Then, as $N \to \infty$*

$$\Pr\{u_t \in B_{U_t}|\{x_{t:T}(i), u_{t:T-1}(i)\}_{i=1}^{N}\} \overset{a.s.}{\to} Q^*_{U_t|X_t}(B_{U_t}|x_t).$$

*Proof.* See Appendix D ∎

We showed that under Assumption 1, the optimal deceptive policies can be synthesized using path integral control. Algorithm 1 allows the deceptive agent to numerically compute optimal actions via Monte Carlo simulations without explicitly synthesizing the policy. Since Monte Carlo simulations can be efficiently parallelized, the agent can generate the optimal control actions online.

## IV. NUMERICAL EXAMPLE

In this section, we validate the path-integral-based algorithm proposed to generate optimal deceptive control actions. The problem is illustrated in Figure 2. A supervisor wants

(a) Paths under $R$, $\mathrm{Pr}^{\mathrm{safe}} = 0.04$

(b) Paths under $\widehat{Q}^*$ with $\lambda = 3$, $\mathrm{Pr}^{\mathrm{safe}} = 0.48$

(c) Paths under $\widehat{Q}^*$ with $\lambda = 2$, $\mathrm{Pr}^{\mathrm{safe}} = 0.62$

(d) Paths under $\widehat{Q}^*$ with $\lambda = 0.5$, $\mathrm{Pr}^{\mathrm{safe}} = 0.94$
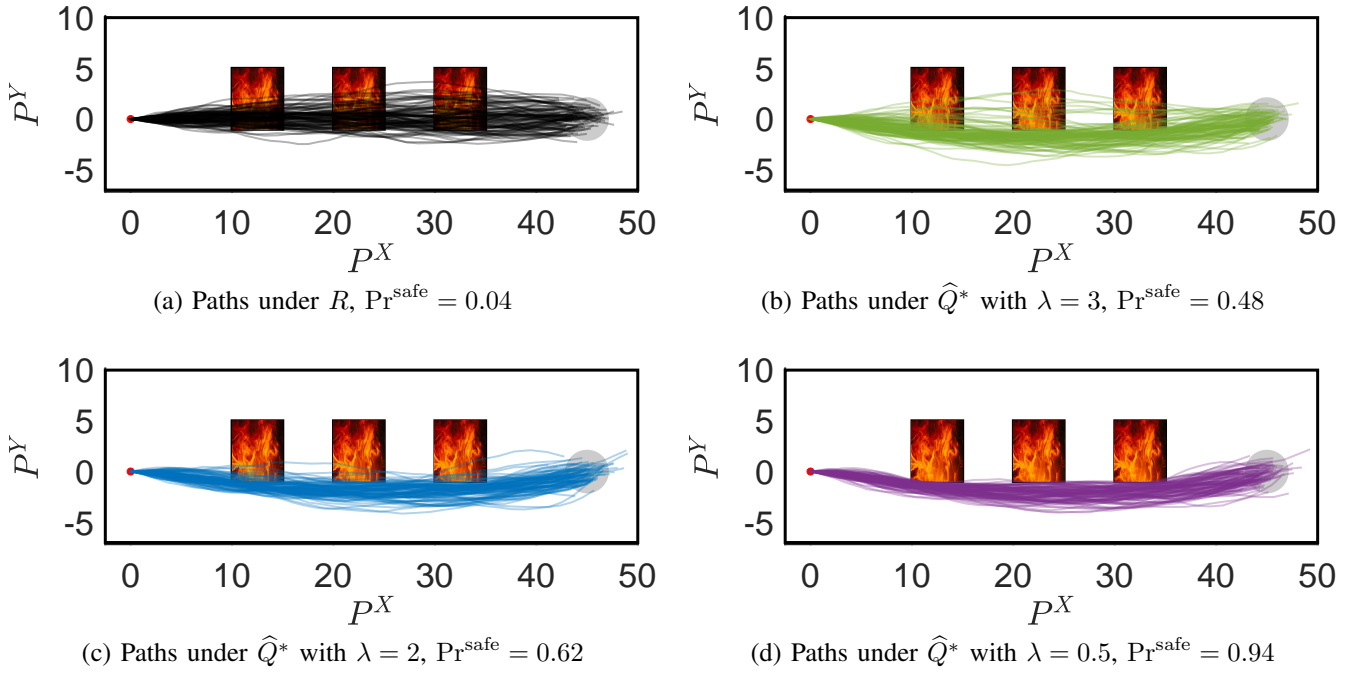
Fig. 2. A unicycle navigation problem. The start position is shown by a red dot, and the goal region by a disk colored in gray. 100 sample paths generated under the reference policy $R$ and the agent's policy $\widehat{Q}^*$ with three values of $\lambda$ are shown. The probability of safe paths $\mathrm{Pr}^{\mathrm{safe}}$ are noted below each case.

an agent to start from the origin and reach a disk of radius $G^R$ centered at $\begin{bmatrix} G^X & G^Y \end{bmatrix}^\top$ (shown in gray color) as fast as possible. The supervisor also expects the agent to inspect the region on the way. To encourage exploration and to provide robustness against unmodeled dynamics, the supervisor designs a randomized reference policy. The agent, on the other hand, wishes to avoid the regions on the way that are covered under fire, as shown in Figure 2. Let these regions be represented collectively by $\mathcal{X}^{\mathrm{fire}}$. Suppose the agent's dynamics are modeled by a unicycle model as:

$$
\begin{aligned}
P_{t+1}^X &= P_{t+1}^X + S_t \cos \Theta_t h \\
P_{t+1}^Y &= P_{t+1}^Y + S_t \sin \Theta_t h \\
S_{t+1} &= S_t + A_t h \\
\Theta_{t+1} &= \Theta_t + \Omega_t h
\end{aligned}
$$

where $(P_t^X, P_t^Y)$, $S_t$, and $\Theta_t$ denote the $x-y$ position, and the heading angle of the agent at time step $t$, respectively. The control input $U_t := \begin{bmatrix} A_t & \Omega_t \end{bmatrix}^T$ consists of acceleration $A_t$ and angular speed $\Omega_t$. $h$ is the time discretization parameter used for discretizing the continuous-time unicycle model. For this simulation study, we set $h = 1$. Note that the agent's dynamics is deterministic as per Assumption 1; however, the control input $U_t$ can be stochastic. Suppose the supervisor designs the reference policy $R$ as a Gaussian probability density with mean $\bar{u}_t$ and covariance $\Sigma_t$:

$$
R_{U_t|X_t}(\cdot|x_t) = \frac{\exp\left[-\frac{1}{2}(u_t - \bar{u}_t)^\top \Sigma_t^{-1}(u_t - \bar{u}_t)\right]}{\sqrt{(2\pi)^2|\Sigma_t|}}.
$$

The mean $\bar{u}_t := \begin{bmatrix} \bar{a}_t & \bar{\omega}_t \end{bmatrix}^\top$ is designed using a proportional controller as

$$
\overline{A}_t = -k_A(S_t - S_t^{\mathrm{desired}}), \quad \overline{\Omega}_t = -k_\Omega(\Theta_t - \Theta_t^{\mathrm{desired}})
$$

where $k_A$ and $k_\Omega$ are proportional gains and $S_t^{\mathrm{desired}}$, $\Theta_t^{\mathrm{desired}}$ are computed as

$$
S_t^{\mathrm{desired}} = \frac{\left\| \begin{bmatrix} G^X \\ G^Y \end{bmatrix} - \begin{bmatrix} P_t^X \\ P_t^Y \end{bmatrix} \right\|}{T - t}, \quad \Theta_t^{\mathrm{desired}} = \tan^{-1}\left(\frac{G^Y - P_t^Y}{G^X - P_t^X}\right).
$$

As mentioned before, the agent wishes to avoid the region $\mathcal{X}^{\mathrm{fire}}$. Suppose the cost function $C_{0:T}$ is designed as

$$
C_{0:T}(X_{0:T}, U_{0:T-1}) = \sum_{t=0}^{T} \mathbb{1}_{[P_t^X \ P_t^Y]^\top \in \mathcal{X}^{\mathrm{fire}}}
$$

where $\mathbb{1}_{[P_t^X \ P_t^Y]^\top \in \mathcal{X}^{\mathrm{fire}}}$ represents an indicator function that returns 1 when the agent is inside the region $\mathcal{X}^{\mathrm{fire}}$ and 0 otherwise. For this simulation, we set

$$
\begin{bmatrix} G^X \\ G^Y \end{bmatrix} = \begin{bmatrix} 45 \\ 0 \end{bmatrix}, \quad \Sigma_t = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}, \quad k_A = 0.1, \ k_\Omega = 0.2, \ T = 50.
$$

The agent chooses its action at each time step using Algorithm 1 where the number of samples is $N = 10^5$.

Suppose $\widehat{Q}^*$ denotes the deceptive agent's distribution generated by the sampling-based Algorithm 1. Figure 2 shows 100 paths under the reference distribution R (Figure 2(a)) and the agent's distribution $\widehat{Q}^*$ for three values of $\lambda$ (Figure 2(b) - 2(d)). A lower value of $\lambda$ implies that the agent cares less about its deviation from the reference policy and more about avoiding the region $\mathcal{X}^{\mathrm{fire}}$. A higher value of $\lambda$ implies the opposite. We also report $\mathrm{Pr}^{\mathrm{safe}}$, the percentage of paths that avoid $\mathcal{X}^{\mathrm{fire}}$. Under the reference distribution $R$, only 4% of the paths are safe. On the other hand, more paths are safe under the agent's distribution $\widehat{Q}^*$, and as the value of $\lambda$ reduces, $\mathrm{Pr}^{\mathrm{safe}}$ increases.
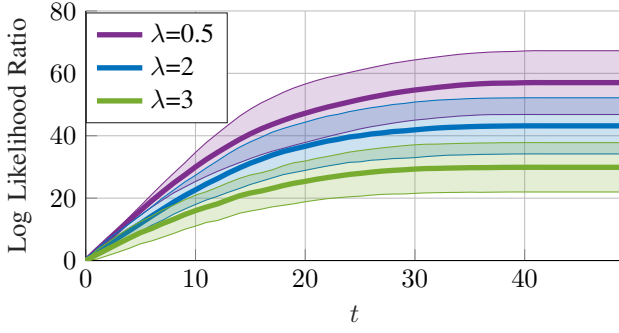
Fig. 3. Expected LLR (with one standard deviation) with respect to time $t$ for three values of $\lambda$.

Figure 3 shows the expected log-likelihood ratio (with one standard deviation) with respect to time $t$ for three values of $\lambda$. The expected LLR is computed as follows. Algorithm 1 selects a control input $u_k \leftarrow u_k(j_k)$ at time step $k$, where $j_k$ is a sample ID obtained from step 5. From the construction of the algorithm, at each time step $k$, the probability of choosing the control input $u_k \leftarrow u_k(j_k)$ under the agent's distribution $\widehat{Q}^*$ is $r_k(j_k)/r_k$, where $r_k(j_k)$ and $r_k$ are computed by steps 3 and 4 of Algorithm 1. Whereas the probability of choosing the control input $u_k \leftarrow u_k(j_k)$ under the reference distribution is $1/N$. Therefore, using (22), the expected LLR upto time $t \in \mathcal{T}$ can be approximately computed as

$$\mathbb{E}_{Q^*}\left[\log \frac{dQ^*_{X_{0:t} \times U_{0:t-1}}}{dR_{X_{0:t} \times U_{0:t-1}}}(x_{0:t}, u_{0:t-1})\right]$$

$$=\mathbb{E}_{Q^*}\left[\sum_{k=0}^{t-1} \log \frac{dQ^*_{U_k|X_k}}{dR_{U_k|X_k}}(x_k, u_k)\right] \approx \frac{1}{N_{\widehat{Q}^*}}\sum_{i=1}^{N_{\widehat{Q}^*}}\sum_{k=0}^{t-1}\frac{r_k(j_k)/r_k}{1/N}.$$

where $N_{\widehat{Q}^*}$ is the number of paths generated by repeatedly running Algorithm 1. Note that since we assume the system dynamics to be deterministic (Assumption 1), once the control input $u_k$ is chosen at time step $k$, the state $x_{k+1}$ is uniquely determined. Therefore, while computing the expected LLR, we only need to consider the probabilities of choosing the control input $u_k$ under policies $\widehat{Q}^*$ and $R$. Figure 3 shows that for a lower value of $\lambda$, the expected LLR is higher, i.e., more deviation of $\widehat{Q}^*$ from $R$.

## V. CONCLUSION

We presented a deception problem under supervisory control for continuous-state discrete-time stochastic systems. Using motivations from hypothesis testing theory, we formalized the synthesis of an optimal deceptive policy as a KL control problem and solved it using backward dynamic programming. Since the dynamic programming approach suffers from the curse of dimensionality, we proposed a simulator-driven algorithm to compute optimal deceptive actions via path integral control. The proposed approach allows the agent to numerically compute deceptive actions online via Monte Carlo sampling of system paths. We validated the proposed approach via a numerical example with a nonlinear system.

For future work, we plan to study the deception problem for continuous-time stochastic systems. We also plan to conduct the sample complexity analysis of the path integral approach to solve KL control problems.

## APPENDIX

### A. Proof of Equation (6)

$$\int_B Q_{X_{0:T} \times U_{0:T-1}}(dx_{0:T} \times du_{0:T-1})$$

$$= \int_B \prod_{t=0}^{T-1} P(dx_{t+1}|x_t, u_t)Q(du_t|x_t) \tag{20a}$$

$$= \int_B \left(\prod_{t=0}^{T-1}\frac{dQ_{U_t|X_t}}{dR_{U_t|X_t}}(x_t, u_t)\right)\prod_{t=0}^{T-1}P(dx_{t+1}|x_t, u_t)R(du_t|x_t) \tag{20b}$$

$$= \int_B \left(\prod_{t=0}^{T-1}\frac{dQ_{U_t|X_t}}{dR_{U_t|X_t}}(x_t, u_t)\right)R_{X_{0:T} \times U_{0:T-1}}(dx_{0:T} \times du_{0:T-1}) \tag{20c}$$

where, $B$ is a Borel set belonging to the $\sigma-$algebra $\mathcal{B}(\mathcal{X}_{0:T} \times \mathcal{U}_{0:T-1})$. The first equality (20a) follows from the definition (3), the second equality (20b) by the definition of the Radon-Nikodym derivative [21] and the last one (20c) from the definition (2). Using (20), we can write the Radon-Nikodym derivative $\frac{dQ_{X_{0:T} \times U_{0:T-1}}}{dR_{X_{0:T} \times U_{0:T-1}}}$ as follows:

$$\frac{dQ_{X_{0:T} \times U_{0:T-1}}}{dR_{X_{0:T} \times U_{0:T-1}}}(x_{0:T}, u_{0:T-1}) = \prod_{t=0}^{T-1}\frac{dQ_{U_t|X_t}}{dR_{U_t|X_t}}(x_t, u_t). \tag{21}$$

Using (21), we get the following:

$$D(Q\|P) = \mathbb{E}_Q\left[\log \frac{dQ_{X_{0:T} \times U_{0:T-1}}}{dR_{X_{0:T} \times U_{0:T-1}}}(x_{0:T}, u_{0:T-1})\right]$$

$$= \mathbb{E}_Q\left[\log \prod_{t=0}^{T-1}\frac{dQ_{U_t|X_t}}{dR_{U_t|X_t}}(x_t, u_t)\right]$$

$$= \mathbb{E}_Q\left[\sum_{t=0}^{T-1}\log \frac{dQ_{U_t|X_t}}{dR_{U_t|X_t}}(x_t, u_t)\right] \tag{22}$$

$$= \mathbb{E}_Q\left[\sum_{t=0}^{T-1}D(Q_{U_t|X_t}(\cdot|X_t)\|R_{U_t|X_t}(\cdot|X_t))\right].$$

### B. Legendre Duality

Let $P$ and $Q$ be probability distributions on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$, and $C : \mathcal{X} \to \mathbb{R}$ a given cost function. Define the internal energy $U(P, C)$, free energy $F(R, C)$ and relative entropy (KL divergence) $D(P\|R)$ as:

$$U(P, C) := \int_{\mathcal{X}} C(x)P(dx)$$

$$F(R, C) := -\lambda \log \int_{\mathcal{X}} \exp\left(-\frac{C(x)}{\lambda}\right)R(dx)$$

$$D(P\|R) := \int_{\mathcal{X}} \log \frac{dP}{dR}(x)P(dx).$$

Then the following duality relationship holds:

$$F(R, C) = \inf_P \{U(P, C) + \lambda D(P\|R)\}$$

$$-\lambda D(P\|R) = \inf_C \{U(P, C) - F(R, C)\}.$$

Also, the optimal probability distribution $P^*$ is given by

$$P^*(B) = \frac{\int_B \exp(-C(x)/\lambda)R(dx)}{\int_{\mathcal{X}} \exp(-C(x)/\lambda)R(dx)}, \quad \forall B \in \mathcal{B}(\mathcal{X}).$$

See [22], [23] for further discussions.

### C. Proof of Theorem 1

By Bellman's optimality principle, the value function satisfies the following recursive relationship:

$$J_t(x_t) = \inf_{Q_{U_t|X_t}} \int_{\mathcal{U}_t} \left\{ \rho_t(x_t, u_t) + \lambda \log \frac{dQ}{dR}(u_t|x_t) \right\} Q(du_t|x_t) \tag{23}$$

where $\rho_t(x_t, u_t)$ is defined by (12). Invoking the Legendre duality between the KL divergence and free energy (see Appendix B), it can be shown that there exists a minimizer $Q^*_{U_t|X_t}$ of the right-hand side of (23), which can be written as

$$Q^*_{U_t|X_t}(B_{U_t}|x_t) = \frac{\int_{B_{U_t}} \exp(-\rho_t(x_t, u_t)/\lambda)R(du_t|x_t)}{\int_{\mathcal{U}_t} \exp(-\rho_t(x_t, u_t)/\lambda)R(du_t|x_t)} \tag{24}$$

where $B_{U_t}$ is a Borel set belonging to the $\sigma-$algebra $\mathcal{B}(\mathcal{U}_t)$. Using (24), the value of (23) can be computed as

$$J_t(x_t) = -\lambda \log \left\{ \int_{\mathcal{U}_t} \exp\left(-\frac{\rho_t(x_t, u_t)}{\lambda}\right) R(du_t|x_t) \right\}. \tag{25}$$

Substituting (12) into (25), we obtain the recursive expression (10).

### D. Proof of Theorem 2

Let $\mathcal{I}_{B_{U_t}}$ be the set of indices of sample paths for which an action in $B_{U_t}$ is taken at time step $t$, i.e., $\mathcal{I}_{B_{U_t}} = \{i \in \{1, 2, \ldots, N\} | u_t(i) \in B_{U_t}\}$. Define the sum of the exponentiated path costs of the paths in $\mathcal{I}_{B_{U_t}}$ as $r_{B_{U_t}} = \sum_{i \in \mathcal{I}_{B_{U_t}}} r_t(i)$. By construction of Algorithm 1,

$$\Pr\{u_t \in B_{U_t} | \{x_{t:T}(i), u_{t:T-1}(i)\}_{i=1}^N\} = \frac{r_{B_{U_t}}}{r_t}. \tag{26}$$

Now, from (17), as $N \to \infty$, we get

$$\frac{r_t}{N} = \frac{1}{N} \sum_{i=1}^N \exp\left(-\frac{C_{t:T}(x_{t:T}(i), u_{t:T-1}(i))}{\lambda}\right) \overset{a.s.}{\to} Z_t(x_t).$$

Similarly, as $N \to \infty$,

$$\frac{r_{B_{U_t}}}{N} = \frac{1}{N} \sum_{i \in \mathcal{I}_{B_{U_t}}} \exp\left(-\frac{C_{t:T}(x_{t:T}(i), u_{t:T-1}(i))}{\lambda}\right)$$

$$\overset{a.s.}{\to} \int_{\{\mathcal{X}_{t+1:T}, \mathcal{U}_{t:T-1}|u_t \in B_{U_t}\}} \exp\left(-\frac{C_{t:T}(x_{t:T}, u_{t:T-1})}{\lambda}\right)$$
$$\times R(dx_{t+1:T}, du_{t:T-1}|x_t)$$

Thus, from (18c) and (26)

$$\Pr\{u_t \in B_{U_t} | \{x_{t:T}(i), u_{t:T-1}(i)\}_{i=1}^N\}$$

$$\overset{a.s.}{\to} \frac{1}{Z_t(x_t)} \int_{\{\mathcal{X}_{t+1:T}, \mathcal{U}_{t:T-1}|u_t \in B_{U_t}\}} \exp\left(-\frac{C_{t:T}(x_{t:T}, u_{t:T-1})}{\lambda}\right)$$
$$\times R(dx_{t+1:T}, du_{t:T-1}|x_t)$$

$$= Q^*_{U_t|X_t}(B_{U_t}|x_t).$$

## REFERENCES

[1] T. M. Cover, *Elements of information theory*. John Wiley & Sons, 1999.

[2] J. Bretagnolle and C. Huber, "Estimation des densités: risque minimax," *Séminaire de probabilités de Strasbourg*, vol. 12, pp. 342–363, 1978.

[3] J. Shim and R. C. Arkin, "A taxonomy of robot deception and its benefits in HRI," in *2013 IEEE international conference on systems, man, and cybernetics*. IEEE, 2013, pp. 2328–2335.

[4] A. Dragan, R. Holladay, and S. Srinivasa, "Deceptive robot motion: synthesis, analysis and experiments," *Autonomous Robots*, vol. 39, pp. 331–345, 2015.

[5] M. O. Karabag, M. Ornik, and U. Topcu, "Deception in supervisory control," *IEEE Transactions on Automatic Control*, vol. 67, no. 2, pp. 738–753, 2021.

[6] ——, "Exploiting partial observability for optimal deception," *IEEE Transactions on Automatic Control*, 2022.

[7] M. Lloyd, *The art of military deception*. Pen and Sword, 2003.

[8] C. Wang and Z. Lu, "Cyber deception: Overview and the road ahead," *IEEE Security & Privacy*, vol. 16, no. 2, pp. 80–85, 2018.

[9] C. Keroglou and C. N. Hadjicostis, "Probabilistic system opacity in discrete event systems," *Discrete Event Dynamic Systems*, vol. 28, pp. 289–314, 2018.

[10] E. Kung, S. Dey, and L. Shi, "The performance and limitations of epsilon-stealthy attacks on higher order systems," *IEEE Transactions on Automatic Control*, vol. 62, no. 2, pp. 941–947, 2016.

[11] C.-Z. Bai, F. Pasqualetti, and V. Gupta, "Data-injection attacks in stochastic control systems: Detectability and performance tradeoffs," *Automatica*, vol. 82, pp. 251–260, 2017.

[12] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International conference on machine learning*. PMLR, 2015, pp. 1889–1897.

[13] S. Filippi, O. Cappé, and A. Garivier, "Optimism in reinforcement learning and Kullback-Leibler divergence," in *2010 48th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, 2010, pp. 115–122.

[14] E. Todorov, "Linearly-solvable Markov decision problems," *Advances in neural information processing systems*, pp. 1369–1376, 2007.

[15] K. Ito and K. Kashima, "Kullback–Leibler control for discrete-time nonlinear systems on continuous spaces," *SICE Journal of Control, Measurement, and System Integration*, vol. 15, no. 2, pp. 119–129, 2022.

[16] H. J. Kappen, "Path integrals and symmetry breaking for optimal control theory," *Journal of statistical mechanics: theory and experiment*, vol. 2005, no. 11, p. P11011, 2005.

[17] E. Theodorou, J. Buchli, and S. Schaal, "A generalized path integral control approach to reinforcement learning," *The Journal of Machine Learning Research*, vol. 11, pp. 3137–3181, 2010.

[18] G. Williams, P. Drews, B. Goldfain, J. M. Rehg, and E. A. Theodorou, "Aggressive driving with model predictive path integral control," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 1433–1440.

[19] G. Williams, A. Aldrich, and E. A. Theodorou, "Model predictive path integral control: From theory to parallel computation," *Journal of Guidance, Control, and Dynamics*, vol. 40, no. 2, pp. 344–357, 2017.

[20] D. Bertsekas, *Dynamic programming and optimal control: Volume I*. Athena scientific, 2012, vol. 1.

[21] R. Durrett, *Probability: theory and examples*. Cambridge university press, 2019, vol. 49.

[22] M. Boué and P. Dupuis, "A variational representation for certain functionals of Brownian motion," *The Annals of Probability*, vol. 26, no. 4, pp. 1641–1659, 1998.

[23] E. A. Theodorou and E. Todorov, "Relative entropy and free energy dualities: Connections to path integral and KL control," *The 51st IEEE Conference on Decision and Control (CDC)*, pp. 1466–1473, 2012.