



Stochastic control with distributionally robust constraints for cyber–physical systems vulnerable to attacks[☆]

Nishanth Venkatesh^{a,*}, Aditya Dave^{b,2}, Ioannis Faros^{a,1}, Andreas A. Malikopoulos^{a,b,3}

^a Department of Systems Engineering, Cornell University, Ithaca, NY 14850, USA

^b School of Civil and Environmental Engineering, Cornell University, Ithaca, NY 14853, USA

ARTICLE INFO

Recommended by T. Parisini

Keywords:

Stochastic control
Constrained MDPs
Distributionally robust control
Cyber–physical systems

ABSTRACT

In this paper, we investigate the control of a cyber–physical system (CPS) while accounting for its vulnerability to external attacks. We formulate a constrained stochastic problem with a robust constraint to ensure robust operation against potential attacks. We seek to minimize the expected cost subject to a constraint limiting the worst-case expected damage an attacker can impose on the CPS. We present a dynamic programming decomposition to compute the optimal control strategy in this robust-constrained formulation and prove its recursive feasibility. We also illustrate the utility of our results by applying them to a numerical simulation.

1. Introduction

Cyber–physical systems (CPSs) have enabled highly efficient control of physical processes by tightly coupling sensing, communication, and computational processing to generate real-time decisions with classical (Kim & Kumar, 2012) and nonclassical information structures (Malikopoulos, 2023a). They span various important applications including, but not limited to, connected and automated vehicles (Malikopoulos, Beaver, & Chremos, 2021; Venkatesh, Le, Dave, & Malikopoulos, 2023), Internet of Things (Ansere, Han, Liu, Peng, & Kamal, 2020), and social media platforms (Dave, Chremos, & Malikopoulos, 2022). However, in each of these applications, the interplay between the cyber components and the physical world can make the system vulnerable to various security threats, e.g., control system malware (Baezner & Robin, 2017) and staged attacks (Serror, Hack, Henze, Schuba, & Wehrle, 2020). This has led to many studies on controlling CPSs while ensuring robustness and resilience to attacks (Ghiasi et al., 2023; Zhang et al., 2023).

The common modeling framework for CPSs utilizes a stochastic formulation to account for uncertainties in the dynamics that arise within the evolution of the physical process. In this formulation, an agent is assumed to have access to a prior distribution for all uncertainties and must compute a control strategy to generate real-time control actions that minimize the total expected cost (Dave, Venkatesh,

& Malikopoulos, 2022b; Sutton & Barto, 2018). In stochastic formulations, constraints on the state and actions are modeled as probabilistic constraints, which can be imposed either in expectation or with some probabilities (Varagapriya & Singh, 2023). Similarly, approaches like those reported in Altman (2021), Ermon, Gomes, Selman, and Vladimirovsky (2012) consider probabilistic constraints on the cumulative reward. However, the actual performance and constraint satisfaction of an optimal control strategy are very sensitive to changes in a mismatch between the assumed prior probability model and the actual model (Malikopoulos, 2023b; Mannor, Simester, Sun, & Tsitsiklis, 2007; Wiesemann, Kuhn, & Rustem, 2013). Such a mismatch is bound to occur when a CPS is under attack from an adversary. Thus, it may not be appropriate to model safety–critical requirements on system behavior using probabilistic constraints in a stochastic formulation.

To accommodate the needs of safety–critical systems, several research efforts (Bertsekas & Rhodes, 1973; Gagrani & Nayyar, 2017; Iyengar, 2005; Shoukry, Araujo, Tabuada, Srivastava, & Johansson, 2013) have explored minimax formulations. Similar approaches (Dave, Venkatesh, & Malikopoulos, 2022c, 2023) consider non-stochastic formulations in which the agent does not have knowledge about the distributions of uncertainties and uses only the set of feasible values to compute optimal strategies that minimize the maximum costs. Though such approaches are suitable for applications under attack, such as cyber-security (Rasouli, Miehl, & Teneketzis, 2018), and power systems (Zhu & Başar, 2011), during regular operation of systems without

[☆] This research was supported by NSF under Grants CNS-2149520 and CMMI-2219761.

* Corresponding author.

E-mail addresses: ns942@cornell.edu (N. Venkatesh), a.dave@cornell.edu (A. Dave), if74@cornell.edu (I. Faros), amaliko@cornell.edu (A.A. Malikopoulos).

¹ Student Member, IEEE.

² Member, IEEE.

³ Senior Member, IEEE.

attacks, they lead to outcomes that are overly conservative (Coraluppi & Marcus, 2000). Consequently, there remains a need for alternative approaches towards the control of vulnerable systems that avoid overly conservative decision-making during regular system operation and maintain a level of reliability when the system is occasionally attacked.

In this paper, we combine the superior performance of stochastic formulations in achieving an objective and the safety guarantees of worst-case formulations (Dave, Venkatesh, & Malikopoulos, 2022a) in minimizing vulnerabilities. To this end, we impose a distributionally robust constraint on a secondary objective that accounts for the vulnerabilities of a CPS to an attack. Concurrently, we aim to minimize the expected value of a primary cost for the best performance over a finite horizon. Our formulation generalizes the previous work of Chen and Blankenship (2004), which addressed the problem of minimizing an expected discounted cost subject to either an expected or a minimax constraint. We consider an uncertainty set which is the set of possible probability distributions (Clement & Kroer, 2021) from which an attack could happen. By considering a distributionally robust constraint, our formulation allows for greater control over the trade-off between conservativeness and optimality by appropriately adjusting the size of the uncertainty set for probability distributions. In the extreme case that the set of feasible distributions is a singleton, we recover an expected value constraint. In contrast, if we expand the set to allow every possible distribution on the state space, we recover the non-stochastic worst-case constraint as a special case. Thus, by changing the set of feasible distributions, we can better select the level of conservativeness of our formulation.

Our main contributions in this paper are (1) the problem formulation of controlling a vulnerable CPS using a stochastic cost and distributionally robust constraint (Problem 1), (2) a dynamic programming (DP) decomposition for this problem, which computes the optimal strategy that ensures recursive feasibility of the constraint (Theorem 1), and (3) the illustration of the utility of our results by comparing them to both stochastic and worst-case approaches in numerical simulation (Section 4).

The remainder of the paper proceeds as follows. In Section 2, we formulate the problem. In Section 3, we present the DP decomposition. In Section 4, we demonstrate our results in a numerical example, and in Section 5, we draw concluding remarks.

2. Model

We consider a CPS whose evolution is described by a finite Markov decision process (MDP), denoted by a tuple $(\mathcal{X}, \mathcal{U}, n, P, c, c_n)$, where \mathcal{X} is a finite state space and \mathcal{U} is a finite set of feasible actions available to an agent seeking to control the MDP. The system evolves over discrete time steps denoted by $t = 0, \dots, n$, where $n \in \mathbb{N}$ is the finite time horizon. The state of the system and the control action of the agent at each t are denoted by the random variables X_t and U_t , respectively. The transition function at each t is denoted by $P_t : \mathcal{X} \times \mathcal{U} \times \mathcal{X} \rightarrow \Delta(\mathcal{X})$, where $\Delta(\mathcal{X})$ is the set of all probability distributions on the state space \mathcal{X} . Nominally, the transition function is given by $P_t = \bar{P}$ for all t , where $\bar{P} \in \Delta(\mathcal{X})$. For the realizations $x_t \in \mathcal{X}$ and $u_t \in \mathcal{U}$ of the state X_t and the control action U_t , the probability of transitioning to a state $x_{t+1} \in \mathcal{X}$ is $\mathbb{P}(X_{t+1} = x_{t+1} \mid x_t, u_t) = P_t(x_{t+1} \mid x_t, u_t)$. The agent selects the action using a control law $g_t : \mathcal{X} \rightarrow \mathcal{U}$ as $U_t = g_t(X_t)$, where g_t is chosen from the feasible set of control laws at time t , denoted as \mathcal{G}_t . The tuple of control laws denotes the control strategy of the agent $g := (g_0, \dots, g_{n-1})$, where $g \in \mathcal{G}$ and $\mathcal{G} = \prod_{t=0}^{n-1} \mathcal{G}_t$. After selecting the action at each $t = 0, \dots, n-1$, the agent incurs a cost $c(X_t, U_t)$ generated using the function $c : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$. Then, the performance of a strategy g is measured by the total expected cost beginning at an initial state $x_0 \in \mathcal{X}$:

$$J_0(g; x_0) = \mathbb{E}^g \left[\sum_{t=0}^{n-1} c(X_t, U_t) + c_n(X_n) \mid x_0 \right], \quad (1)$$

where $c_n : \mathcal{X} \rightarrow \mathbb{R}$ is the terminal cost, and \mathbb{E}^g denotes the expectation on all the random variables with respect to the probability distributions generated by the choice of control strategy g .

In the context of a CPS, the conventional approach of selecting a control strategy g to minimize the total expected cost (1) may not be adequate to ensure smooth operation, particularly when the CPS is vulnerable to attacks by an adversary. We consider that the presence or absence of an adversary during the system's operation is determined at the onset; however, this information is unknown to the agent. The adversary's influence on the system's dynamics results in a change in the transition probability at each $t = 0, \dots, n-1$ from a known set $\mathcal{P} \subseteq \Delta(\mathcal{X})$. Thus, an attack may be reflected by the choice of the worst transition function from \mathcal{P} . We allow the adversary to attack the system with access to the realization of the state $x_t \in \mathcal{X}$ and action $u_t \in \mathcal{U}$. Note that the nominal transition function \bar{P} belongs to the set \mathcal{P} to allow for the case of no attack.

An agent that observes the presence or absence of an adversary can select either a purely robust or risk-neutral formulation, depending on the current situation. However, a risk-neutral formulation may involve an arbitrarily large risk for the agent and leave the CPS vulnerable during an attack. In contrast, a robust formulation may be too conservative for the majority of situations where no attack occurs. Thus, we impose a robust constraint to *limit* the worst-case damage possible during an attack while minimizing the expected total cost. To this end, the agent incurs a constraint penalty $d(X_t, U_t) \in \mathbb{R}$ at each $t = 0, \dots, n-1$. The total expected worst-case penalty is given by

$$\mathcal{L}_0(g; x_0) = \max_{P_{0:n-1} \in \mathcal{P}^n} \mathbb{E}_{P_{0:n-1}}^g \left[\sum_{t=0}^{n-1} d(X_t, U_t) + d_n(X_n) \mid x_0 \right], \quad (2)$$

where $d_n : \mathcal{X} \rightarrow \mathbb{R}$ is the terminal penalty, $P_{0:n-1}$ is the collection of transition functions for $t = 0, \dots, n-1$, each taking values in the set \mathcal{P} . Note that this penalty has a distributionally robust form where the attacker may select the worst transition function $P_t \in \mathcal{P}$ at each t . Furthermore, the choice of a particular function at any time t does not limit the functions available to the adversary at time $t+1$ in (2). The distributionally robust constraint is formulated by defining an upper bound $l_0 \in \mathbb{R}$, on the worst-case total expected penalty.

Remark 1. During an attack, the agent may prioritize a different property, e.g., safety, of the system rather than the total expected cost used in (1). Hence, the constraint penalty at each instance of time is considered to be distinct from the cost. However, if we seek to limit the influence of the adversary on the performance itself, the penalty in the constraint can be set equal to the cost at each t .

Next, we define the agent's constrained control problem.

Problem 1. The optimization problem is to compute the optimal control strategy $g^* \in \mathcal{G}$, if one exists, subject to a constraint on (2), i.e.,

$$\min_{g \in \mathcal{G}} J_0(g; x_0), \quad (3)$$

$$\text{s.t. } \mathcal{L}_0(g; x_0) \leq l_0, \quad (4)$$

for a given MDP $(\mathcal{X}, \mathcal{U}, n, P, c, c_n)$, penalty functions (d, d_n) , set of transition functions \mathcal{P} , upper bound $l_0 \in \mathbb{R}$, and initial state $x_0 \in \mathcal{X}$.

We impose the following assumptions on our formulation:

Assumption 1. The costs and penalties at each instance of time are upper bounded by the finite maximum values $c^M \in \mathbb{R}$ and $d^M \in \mathbb{R}$, respectively. They are also lower bounded by the finite minimum values $c^m \in \mathbb{R}$ and $d^m \in \mathbb{R}$, respectively.

Assumption 1 ensures that the expected total cost (1) and robust total penalty (2) are finite for any value of $n \in \mathbb{N}$.

Assumption 2. The bound $l_0 \in \mathbb{R}$ is such that the set $\mathcal{G}_{l_0} := \{g \in \mathcal{G} \mid \mathcal{L}_0(g; x_0) \leq l_0\}$ is not empty.

Assumption 2 ensures that Problem 1 has a feasible solution and, thus, it is well-posed. Our goal is to efficiently compute an optimal solution to Problem 1 without violating the constraint. Next, we present a DP decomposition for the problem.

3. Dynamic programming decomposition

In this section, we present the value functions that constitute a DP decomposition to compute the optimal control strategy g^* for Problem 1. To show that the computed strategy satisfies the distributionally robust constraint (4), we need to prove its recursive feasibility for all $t = 1, \dots, n-1$. To achieve this, in Section 3.1, we define the *penalty-to-go function* to express the application of the constraint only from any time t to the terminal time n . We then construct a set of upper bounds on the penalty-to-go function at any $t = 0, \dots, n-1$, such that these bounds admit a feasible solution, and present a methodology to compute these sets. We also introduce the notion of bound functions, which will be utilized to ensure recursive feasibility. In Section 3.2, we use bound functions within the proposed DP decomposition and prove its optimality.

3.1. Feasible bound for robust constraint

We begin by constructing the *penalty-to-go* function that maps each realization of the state $x_t \in \mathcal{X}$ at any $t = 0, \dots, n-1$ to an expected worst-case penalty to reach n using a sequence of control laws $g_{t:n-1} \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell$. Specifically, this penalty-to-go at each t is

$$\mathcal{L}_t(g_{t:n-1}; x_t) = \max_{P_{t:n-1} \in \mathcal{P}^{n-t}} \mathbb{E}_{P_{t:n-1}}^{g_{t:n-1}} \left[\sum_{\ell=t}^{n-1} d(X_\ell, U_\ell) + d_n(X_n) \mid x_t \right], \quad (5)$$

where the expectation on all the random variables is with respect to the distributions generated by the choice of control laws $g_{t:n-1} \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell$. Importantly, the control laws utilized prior to time t do not influence the penalty-to-go from time t . Additionally, note that the penalty-to-go from $t = 0$ is the total expected penalty in (2). Next, we construct a set of feasible upper bounds on the penalty-to-go function.

Definition 1. For all $t = 1, \dots, n-1$, the *set of feasible upper bounds* for a state $x_t \in \mathcal{X}$ is

$$\Lambda_t(x_t) := \left\{ l_t \in \mathbb{R} \mid \exists g_{t:n-1} \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell, \text{ s.t. } \mathcal{L}_t(g_{t:n-1}; x_t) \leq l_t \right\}, \quad (6)$$

with $\Lambda_n(x_n) := [d_n(x_n), d^M]$ at $t = n$ for each $x_n \in \mathcal{X}$ and $\Lambda_0(x_0) := \{l_0\}$ identically for all $x_0 \in \mathcal{X}$.

In Definition 1, the bound l_t acts only upon the penalty-to-go $\mathcal{L}_t(g_{t:n-1}; x_t)$. Thus, each bound $l_t \in \Lambda_t(x_t)$ ensures feasibility of only the control laws $g_{t:n-1} \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell$ for each $x_t \in \mathcal{X}$ and $t = 0, \dots, n-1$.

Next, to ensure recursive feasibility in our solution approach, our goal is to select a feasible bound on the penalty-to-go for all $t = 0, \dots, n-1$. These bounds should ensure that, starting with l_0 at $t = 0$, there exists at least one feasible sequence of control laws $g_{t:n-1} \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell$. We note that based on Assumption 2, such a sequence exists at $t = 0$.

To this end, we establish the notion of *bound functions* $\lambda_t : \mathcal{X} \rightarrow \mathbb{R}$ at each $t = 0, \dots, n$. The output of the bound function $\lambda_t(x_t)$ is a feasible bound from Definition 1 for all $x_t \in \mathcal{X}$ and all t . Then, for any bound $l_t \in \Lambda_t(x_t)$ and a control action $u_t \in \mathcal{U}$, the set of *recursively consistent bound functions* at time $t+1$ is

$$F_t(x_t, u_t, l_t) = \left\{ \lambda_{t+1} \mid \lambda_{t+1}(x_{t+1}) \in \Lambda_{t+1}(x_{t+1}), \forall x_{t+1} \in \mathcal{X} \text{ and } \right.$$

$$\left. \max_{P_t \in \mathcal{P}} \mathbb{E}_{P_t} [\lambda_{t+1}(X_{t+1}) \mid x_t, u_t] \leq l_t - d(x_t, u_t) \right\}. \quad (7)$$

The inequality in the conditioning of the set in (7) yields the allowable bound at time $t+1$ after considering the “consumption” of the bound l_t by the penalty $d(x_t, u_t)$ incurred at time t . This inequality is imposed upon the maximum expected value of $\lambda_{t+1}(X_{t+1})$ given the state x_t and action u_t to ensure recursive constraint satisfaction. Note that this maximization captures the distributionally robust form of transition functions in (5) and, thus, accounts for the possible influence of the attacker. Thus, given $l_t \in \Lambda_t(x_t)$ at time t , restricting attention to $\lambda_{t+1} \in F_t(x_t, u_t, l_t)$ ensures that any selected bound at time $t+1$ is feasible. Beginning with l_0 at $t = 0$ and applying this property for the set $F_t(x_0, u_0, l_0)$ for all $x_0 \in \mathcal{X}$ and $u_0 \in \mathcal{U}$ ensures recursive feasibility and constraint satisfaction for all $t = 0, \dots, n$. Due to the importance of the sets $\Lambda_t(x_t)$ in (7), it is essential to efficiently compute them before deriving an optimal strategy.

To begin, we observe that for any feasible $l_t \in \Lambda_t(x_t)$, there exists a sequence of control laws $g_{t:n-1}$ that satisfies the constraint $\mathcal{L}_t(g_{t:n-1}; x_t) \leq l_t$ for all $l_t \leq \hat{l}_t \in \mathbb{R}$ and all $x_t \in \mathcal{X}$. Hence, it is sufficient to compute the smallest feasible bound $\lambda_t^m(x_t)$ for each $x_t \in \mathcal{X}$ and note that the set $\Lambda_t(x_t) \subseteq [\lambda_t^m(x_t), \infty)$. At the other extreme, without loss of generality, we can restrict the maxima of $\Lambda_t(x_t)$ to $l_t^M = \min\{l_0, \sum_{i=t}^n d_i^M\}$. This is because including bounds larger than l_t^M does not increase the set of feasible sequences of control laws $g_{0:t-1}$. Thus, the structural form of the set of feasible upper bounds is $\Lambda_t(x_t) = [\lambda_t^m(x_t), l_t^M]$ for all $t = 0, \dots, n-1$.

Next, we present a recursive approach to compute $\lambda_t^m(x_t)$ for all t and complete the construction of $\Lambda_t(x_t)$.

Lemma 1. The lower bound $\lambda_t^m(x_t)$ of the set $\Lambda_t(x_t)$ for all $x_t \in \mathcal{X}$ and $t = 1, \dots, n-1$ is obtained by the following minimization problem

$$\lambda_t^m(x_t) = \min_{u_t \in \mathcal{U}} \left\{ d(x_t, u_t) + \max_{P_t \in \mathcal{P}} \mathbb{E}_{P_t} [\lambda_{t+1}^m(X_{t+1}) \mid x_t, u_t] \right\}. \quad (8)$$

Proof. The smallest feasible bound $\lambda_t^m(x_t)$ belongs to the set $\Lambda_t(x_t)$. From Definition 1, we can see that there exists a sequence $g_{t:n-1}$ for which the penalty-to-go is exactly equal to $\lambda_t^m(x_t)$ and any bound smaller than $\lambda_t^m(x_t)$ is infeasible. Thus we can compute $\lambda_t^m(x_t)$ as

$$\lambda_t^m(x_t) = \min_{g_{t:n-1} \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell} \mathcal{L}_t(g_{t:n-1}; x_t), \quad (9)$$

which becomes an instance of the standard distributionally robust DP problem. The objective is to minimize the penalty-to-go while being distributionally robust against the uncertainty in the transition function. Using the arguments presented in Iyengar (2005, Theorem 2.1) for deriving the optimal objective in such a problem, we can see how Lemma 1 computes the minimum value of $\mathcal{L}_t(g_{t:n-1}; x_t)$ at each t . This shows that, Lemma 1 can be used to recursively compute the smallest feasible bound $\lambda_t^m(x_t)$ for each $x_t \in \mathcal{X}$ and all t . \square

3.2. Dynamic program for Problem 1

In this subsection, before presenting the DP decomposition, we begin by defining the cost-to-go in a manner similar to the penalty-to-go in Section 3.1. For all $t = 0, \dots, n-1$, the cost-to-go from any $x_t \in \mathcal{X}$ is

$$J_t(g_{t:n-1}; x_t) = \mathbb{E}_{g_{t:n-1}} \left[\sum_{\ell=t}^{n-1} c(X_\ell, U_\ell) + c_n(X_n) \mid x_t \right], \quad (10)$$

where $\mathbb{E}_{g_{t:n-1}}$ denotes the expectation on all the random variables with respect to the distributions generated by the nominal transition function \bar{P} and the choice of control laws $g_{t:n-1} \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell$. Note that the cost-to-go at time t is affected only by the sequence of control laws $g_{t:n-1}$ and the cost-to-go at $t = 0$ is equivalent to the performance measure (1) for any strategy g .

Before we can construct a DP decomposition, we recall that [Problem 1](#) also restricts the set of feasible strategies by constraining the penalty-to-go $\mathcal{L}_0(g; x_0)$ with an upper bound l_0 . Thus, as derived in [Section 3.1](#), we need to impose a constraint using the bound function $\lambda_t \in F_{t-1}(x_{t-1}, u_{t-1}, l_{t-1})$ for all $t = 1, \dots, n$ to ensure recursive feasibility in our solution approach. To this end, at each $t = 0, \dots, n-1$, we expand our state-space \mathcal{X} by appending with it a set of possible bounds, \mathbb{R} . Thus, the value functions of our DP decomposition are functions of $(X_t, L_t) \in \mathcal{X} \times \mathbb{R}$, where the random variable $L_t = \lambda_t(X_t)$. The realizations of the random variable L_t are denoted by l_t . Furthermore, at each $x_t \in \mathcal{X}$, the control law $g_t \in \mathcal{G}_t$ at each $t = 0, \dots, n-1$ selects a control action $U_t \in \mathcal{U}$ using the expanded state space as $U_t = g_t(X_t, L_t)$.

Remark 2. We note that expanding the state space from X_t to (X_t, L_t) expands the domain of control laws as compared to the standard Markovian control law for regular MDPs. However, the result in [Section 3.1](#) is still valid for control laws with this larger domain because the functions introduced in [3.1](#) depend only on the realization $x_t \in \mathcal{X}$ of X_t and are independent of the realization $l_t = \lambda_t(x_t)$ of L_t .

For all $t = 0, \dots, n-1$, the value function for all $x_t \in \mathcal{X}$ and $l_t \in \Lambda(x_t)$ corresponding to the sequence of control laws $g_{t:n-1}$ is given by

$$V_t^{g_{t:n-1}}(x_t, l_t) = \begin{cases} \mathcal{J}_t(g_{t:n-1}; x_t) & \text{if } \mathcal{L}_t(g_{t:n-1}; x_t) \leq l_t, \\ \kappa & \text{otherwise,} \end{cases} \quad (11)$$

where $\kappa \in \mathbb{R}$ is a large constant that satisfies $\kappa > n \cdot c^M$ and indicates constraint violation by $g_{t:n-1}$. Eventually, when we minimize over the set of strategies, the presence κ will help us exclude infeasible solutions. At the terminal time n , where no actions are allowed, the value function is simply $V_n(x_n, l_n) = c(x_n)$. Then, the optimal value functions for all $x_t \in \mathcal{X}$, $l_t = \lambda_t(x_t)$ and all $t = 0, \dots, n-1$ are

$$V_t(x_t, l_t) = \min_{g_{t:n-1} \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell} V_t^{g_{t:n-1}}(x_t, l_t). \quad (12)$$

Theorem 1. At each $t = 0, \dots, n-1$, for all $x_t \in \mathcal{X}$ and $l_t = \lambda_t(x_t)$, the optimal value function can be recursively computed using the following DP decomposition:

$$V_t(x_t, l_t) = \min_{\substack{u_t \in \mathcal{U}, \\ \lambda_{t+1} \in F_t(x_t, u_t, l_t)}} \left\{ c(x_t, u_t) + \mathbb{E} \left[V_{t+1}(X_{t+1}, \lambda_{t+1}(X_{t+1})) \mid x_t, u_t \right] \right\}, \quad (13)$$

where, at the terminal time $t = n$, the optimal value function is simply given by $V_n(x_n, l_n) = c(x_n)$.

Proof. We prove that the DP decomposition presented in [Theorem 1](#) computes the optimal value function recursively using mathematical induction. At the terminal time, the value function is given by $V_n(x_n, l_n) = c(x_n)$. Suppose that the optimal value function V_{t+1} at time $t+1$ can be computed according to [\(13\)](#). It is enough to show that [\(13\)](#) can be used to compute $V_t(x_t, l_t)$ at time t . We need first to show that the left-hand side in [\(13\)](#) is lower bounded by the right-hand side and vice-versa. As a result, the left-hand side of [\(13\)](#) will be both upper and lower bounded by the expression on the right-hand side. Hence, we conclude that in [\(13\)](#), the left-hand side is equal to the right-hand side. Details of the mathematical arguments are provided in [Appendix](#). \square

Remark 3. We showed that the DP decomposition presented in [\(13\)](#) computes the optimal value function at each $t = 0, \dots, n-1$. Using [Theorem 1](#), we can compute the sequence of optimal control laws $g_{0:n-1}^* \in \prod_{\ell=0}^{n-1} \mathcal{G}_\ell$ which yields the optimal value function $V_0(x_0, l_0)$ at time $t = 0$.

Remark 4. At any $t = 0, \dots, n-1$, and for all $x_t \in \mathcal{X}$ and a feasible bound l_t , [Theorem 1](#) states that the optimal control action is $u_t^* = g_t^*(x_t, l_t)$, i.e., the minimizing argument in [\(13\)](#). Subsequently, the optimal bound function $\lambda_{t+1}^*(\cdot)$ is computed as a function of the state x_t , bound l_t , and optimal action u_t^* . Since the optimal bound function is computed at the preceding time step, during implementation, λ_t^* is available at the onset of time t and the agent ensures that $l_t = \lambda_t^*(x_t)$. Hence, to solve [Problem 1](#), the control action at all t is selected as $u_t^* = g_t^*(x_t, \lambda_t^*(x_t))$. This shows that the optimal control strategy can be selected using $x_t \in \mathcal{X}$ during implementation.

4. Numerical example

In this Section, we illustrate the efficiency of the effectiveness of our approach using a numerical example. We consider a reach-avoid problem where an agent seeks to navigate to a designated cell in a 4×4 grid world while avoiding a different cell in the grid. At each $t = 0, \dots, n$, the agent's position X_t takes values in the set of grid cells:

$$\mathcal{X} = \{(0,0), (0,1), \dots, (3,2), (3,3)\}. \quad (14)$$

The action U_t denotes the agent's direction of movement and takes values in the set:

$$\mathcal{U} = \{(-1,0), (1,0), (0,0), (0,1), (0,-1)\}. \quad (15)$$

Under normal system operation, the agent has a small chance of movement failure by "slipping". This nominal transition function is modeled by considering that at each t , the agent moves in the direction selected by the action U_t with probability 0.8 and may slip by moving in either the clockwise or anticlockwise direction to U_t with probabilities of 0.1 each. The agent does not slip when selecting the action $(0,0)$, i.e., when deciding not to move. Thus, starting at a randomly selected initial state $x_0 \in \mathcal{X}$, the agent's dynamics for all $t = 0, \dots, n-1$ are

$$\bar{P}(x_{t+1} \mid x_t, u_t) = \begin{cases} 0.8 & \text{if } x_{t+1} = x_t + u_t, \\ 0.1 & \text{if } x_{t+1} = x_t + u_t^{\text{cl}}, \\ 0.1 & \text{if } x_{t+1} = x_t - u_t^{\text{cl}}, \end{cases} \quad (16)$$

where if $u_t = (u_t^1, u_t^2)$, then $u_t^{\text{cl}} = (-u_t^2, u_t^1)$ is the clockwise rotation of u_t . If the agent's position X_t is at one of the four corners of the grid or along the edge of the grid, then the agent may only slip in the available directions and not move off the grid. Thus, if only one direction is available, they move in the direction of the selected action with a probability of 0.8 and move in the available direction with a probability of 0.2.

The goal of the agent is to reach the destination cell $(3,2)$ marked by D, while avoiding a "trap" cell $(2,1)$, marked by X in [Fig. 1](#) and [Fig. 2](#). The agent incurs no interim costs, i.e., for all $t = 0, \dots, n-1$:

$$c(X_t, U_t) = 0. \quad (17)$$

However, at time $n = 10$ which is the end of horizon, the agent incurs a terminal cost given by the distance of the agent to the destination:

$$c_n(X_n) = \eta(X_n, (3,2)), \quad (18)$$

where, $(3,2)$ is the destination and $\eta(\cdot, \cdot)$ denotes the Manhattan distance. For all $t = 0, \dots, n$ the agent incurs a penalty of 1 unit if their position coincides with the trap $(2,1)$:

$$d(X_t, U_t) = \mathbb{I}[X_t = (2,1)]. \quad (19)$$

An adversary, if present, attacks the reliability of the agent's actuator. Thus, under an attack, the probability of slipping may increase. We incorporate vulnerability to attacks by defining, on the tuple of actual movements $(u_t, u_t^{\text{cl}}, -u_t^{\text{cl}})$ for a given action $u_t \in \mathcal{U}$, the set of possible probability distributions:

$$P_{\text{sim}} := \{ (0.7, 0.3, 0), (0.7, 0.2, 0.1), \dots, (0.7, 0, 0.3), \\ (0.8, 0.2, 0), (0.8, 0.1, 0.1), (0.8, 0, 0.2) \}. \quad (20)$$

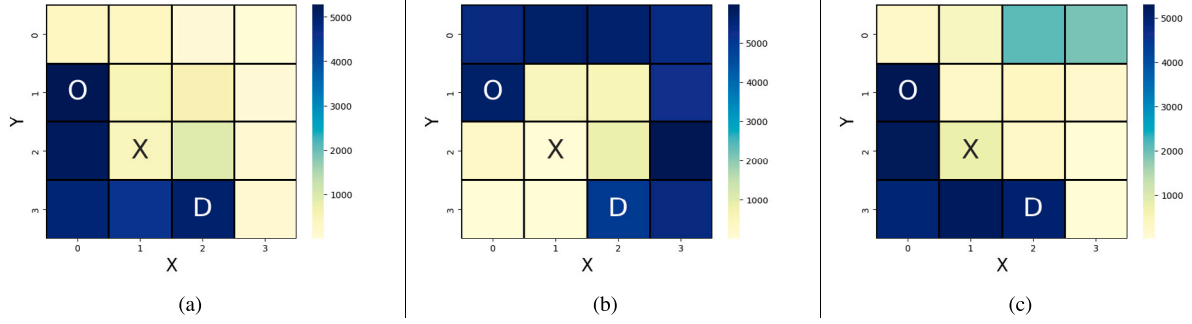


Fig. 1. For the initial state (1,0), the strategy implemented in : (a) distributionally robust (b) conservative (c) stochastic.

Table 1

The number of “trap” cell visits in each case.

Control strategy	Number of visits to “trap” cell	
	$X_0 = (1, 0)$	$X_0 = (0, 1)$
Distributionally Robust: $\mathcal{P} = \mathcal{P}_{sim}$	497	250
Conservative: $\mathcal{P} = \Delta(\mathcal{X})$	35	10
Stochastic: $\mathcal{P} = \{\bar{P}\}$	764	363

We consider two initial positions, (1,0) and (0,1) marked by “O” as shown in Fig. 1 and Fig. 2 respectively. An upper bound of $l_0 = 2.5$ is set on the worst-case total expected penalty. We compare three cases to show how our approach provides more control over the trade-off between conservativeness and optimality for each initial condition. In each of the three cases considered, we specify the set of possible probability distributions \mathcal{P} in (8) to emulate various levels of conservativeness. To demonstrate this, we implement the computed optimal control strategy in a receding horizon manner for 200 time steps. We run 5000 simulations and plot the heat map of the path selected by the agent in Fig. 1 and Fig. 2. In the first case, while we compute the control strategy, we consider $\mathcal{P} = \mathcal{P}_{sim}$, which yields the distributionally robust control strategy. In this case, the agent visits the “trap” cell 497 and 250 times, as shown in Fig. 1(a) and Fig. 2(a), respectively. For the second case, during the computation of the control strategy, we expand the set \mathcal{P} to include every possible distribution in $\Delta(\mathcal{X})$ to yield a conservative strategy. As a result, in Fig. 1(b) and Fig. 2(b), the number of times the agent moves into the “trap” cell are 35 and 10, respectively. Lastly, we compute a stochastic strategy by considering that the set of probability distributions \mathcal{P} is a singleton, with only the nominal transition function. Accordingly, in Fig. 1(c) and Fig. 2(c), the agent moves 764 and 363 times into the “trap” cell, respectively. We present a summary of these results in table 1.

We observe that when the agent utilizes the distributionally robust control strategy, it visits the trap cell more often than the conservative strategy. However, it reaches the target cell in fewer moves than the conservative strategy. On the other hand, it reaches the destination as quickly as the stochastic strategy, with fewer visits to the “trap”, essentially being more robust.

5. Conclusion

In this paper, we proposed the problem of controlling a CPS, which is vulnerable to attack as a distributionally robust stochastic cost minimization problem. For this problem, we presented DP decomposition to compute the optimal control strategy, which ensures the recursive feasibility of the distributionally robust constraint. Finally, we illustrated the utility of our solution approach using a numerical example. Future work should consider using these results in tandem with fast computation techniques for applications with large state space like human-robot collaboration tasks, power grids, and connected and automated vehicles. In such applications, it is essential to avoid over-conservatism while maintaining resilience against any vulnerabilities.

CRedit authorship contribution statement

Nishanth Venkatesh: Conceptualization, Formal analysis, Investigation, Validation, Writing – original draft. **Aditya Dave:** Conceptualization, Formal analysis, Validation, Writing – original draft. **Ioannis Faros:** Software, Visualization. **Andreas A. Malikopoulos:** Conceptualization, Funding acquisition, Project administration, Supervision, Writing – original draft, Writing – review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Appendix. Proof of Theorem 1

We present a detailed proof of Theorem 1 by elaborating on how we show that left-hand side of (13) is both upper and lower bounded by the expression on the right-hand side.

To begin, we recall the induction assumption that the optimal value function V_{t+1} at time $t+1$ can be computed according to (13).

At time t , for each $x_t \in \mathcal{X}$ and a feasible bound $l_t = \lambda_t(x_t)$, let $g_{t:n-1}^* \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell$ be the sequence of control laws such that

$$g_{t:n-1}^* = \arg \min_{g_{t:n-1} \in \prod_{\ell=t}^{n-1} \mathcal{G}_\ell} V_t^{g_{t:n-1}}(x_t, l_t). \quad (21)$$

The value function corresponding to this sequence of control laws is the optimal value function at time t as given by

$$V_t(x_t, l_t) = V_t^{g_{t:n-1}^*}(x_t, l_t). \quad (22)$$

We expand the expression for the value function of this sequence of control laws as

$$\begin{aligned} V_t(x_t, l_t) &= \mathbb{E}^{g_{t:n-1}^*} \left[\sum_{\ell=t}^{n-1} c(X_\ell, U_\ell) + c_n(X_n) \mid x_t \right], \\ &= c(x_t, g_t^*(x_t, l_t)) + \mathbb{E}^{g_{t:n-1}^*} \left[\sum_{\ell=t+1}^{n-1} c(X_\ell, U_\ell) + c_n(X_n) \mid x_t \right], \end{aligned} \quad (23)$$

where, we note that the penalty-to-go of the sequence $g_{t:n-1}^*$ upper bounded by l_t . Hence, we only analyze the cost-to-go component of the value function for this sequence of control laws. We use the law of iterated expectations to introduce X_{t+1} into the inner expectation as

$$\begin{aligned} V_t(x_t, l_t) &= c(x_t, g_t^*(x_t, l_t)) + \\ &\mathbb{E} \left[\mathbb{E}^{g_{t+1:n-1}^*} \left[\sum_{\ell=t+1}^{n-1} c(X_\ell, U_\ell) + c_n(X_n) \mid X_{t+1}, \right. \right. \\ &\left. \left. x_t, g_t^*(x_t, l_t) \right] \right]. \end{aligned} \quad (24)$$

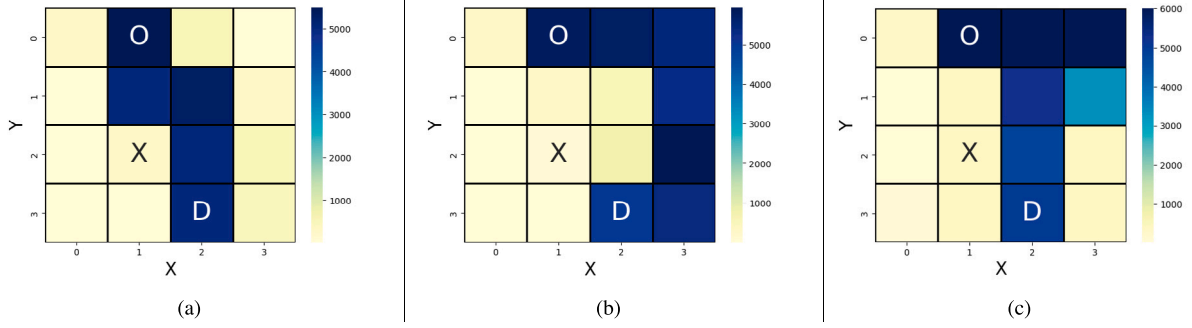


Fig. 2. For the initial state (0, 1), the strategy implemented in : (a) distributionally robust (b) conservative (c) stochastic.

Next, using (10), we write the inner expectation as the cost-to-go from time $t + 1$ for the sequence $g_{t+1:n-1}^*$:

$$V_t(x_t, l_t) = c(x_t, g_t^*(x_t, l_t)) + \mathbb{E} \left[J_{t+1}(g_{t+1:n-1}^*; X_{t+1}) \mid x_t, g_t^*(x_t, l_t) \right]. \quad (25)$$

From (11) and (12), we can see that the cost-to-go for $g_{t+1:n-1}^*$ is always as big as the optimal value function at $t + 1$. Considering this, in addition to the induction hypothesis, we can write:

$$V_t(x_t, l_t) \geq c(x_t, g_t^*(x_t, l_t)) + \mathbb{E} \left[V_{t+1}(X_{t+1}, \mathcal{L}_{t+1}(g_{t+1:n-1}^*; X_{t+1})) \mid x_t, g_t^*(x_t, l_t) \right], \quad (26)$$

where, we note that $\mathcal{L}_{t+1}(g_{t+1:n-1}^*; X_{t+1})$ is exactly the penalty-to-go for the sequence of control laws $g_{t+1:n-1}^*$, hence it is a feasible bound at time $t + 1$. By minimizing the right-hand side of (26) over feasible combination of the action u_t and future bound function λ_{t+1} we have:

$$V_t(x_t, l_t) \geq \min_{\substack{u_t \in U_t \\ \lambda_{t+1} \in F_t(x_t, u_t, l_t)}} \left\{ c(x_t, u_t) + \mathbb{E} \left[V_{t+1}(X_{t+1}, \lambda_{t+1}(X_{t+1})) \mid x_t, u_t \right] \right\}. \quad (27)$$

This shows that the left-hand side of (13), is greater than the right-hand side.

For each time t and $x_t \in \mathcal{X}$ with the bound imposed being $l_t = \lambda_t(x_t)$, we consider

$$(u_t^*, \lambda_{t+1}^*) = \arg \min_{\substack{u_t \in U_t \\ \lambda_{t+1} \in F_t(x_t, u_t, l_t)}} \left\{ c(x_t, u_t) + \mathbb{E} \left[V_{t+1}(X_{t+1}, \lambda_{t+1}(X_{t+1})) \mid x_t, u_t \right] \right\}, \quad (28)$$

where, u_t^* is the optimal action at time t and λ_{t+1}^* is the corresponding optimal bound function at time $t + 1$. Now, we construct a sequence of control laws $\hat{g}_{t:n-1}$ such that:

$$\hat{g}_t(x_t, l_t) = u_t^* \quad (29)$$

$$\hat{g}_{t+\ell} = g_{t+\ell}^*, \text{ for } \ell = 1, \dots, (n-1-t). \quad (30)$$

This implies that at time $t + 1$, the newly constructed sequence of control laws $\hat{g}_{t+1:n-1}$ satisfies the following constraint

$$\mathcal{L}_t(\hat{g}_{t+1:n-1}; x_{t+1}) \leq \lambda_{t+1}^*(x_{t+1}) \forall x_{t+1} \in \mathcal{X}. \quad (31)$$

Next, by using (10), we establish the cost-to-go of the sequence of control laws $\hat{g}_{t:n-1}$ as an upper bound to the optimal value function at t . This is given by

$$V_t(x_t, l_t) \leq J_t(\hat{g}_{t:n-1}; x_t)$$

$$= \mathbb{E}^{\hat{g}_{t:n-1}} \left[\sum_{\ell=t}^{n-1} c(X_\ell, U_\ell) + c_n(X_n) \mid x_t \right]. \quad (32)$$

Then, we use the law of iterated expectations on the right-hand side of (32) to introduce the random variable X_{t+1} , which yields

$$V_t(x_t, l_t) \leq c(x_t, \hat{g}_t(x_t, l_t)) + \mathbb{E} \left[\mathbb{E}^{\hat{g}_{t+1:n-1}} \left[\sum_{\ell=t+1}^{n-1} c(X_\ell, U_\ell) + c_n(X_n) \mid X_{t+1} \right] \mid x_t, \hat{g}_t(x_t, l_t) \right], \quad (33)$$

where we note that the inner expectation only depends on the choice of the sequence $\hat{g}_{t+1:n-1}$. We use (11) and (31) to write the inner expectation as the value function of this sequence of control laws, given by

$$V_t(x_t, l_t) \leq c(x_t, \hat{g}_t(x_t, l_t)) + \mathbb{E} \left[V_{t+1}^{\hat{g}_{t+1:n-1}}(X_{t+1}, \lambda_{t+1}^*(X_{t+1})) \mid x_t, \hat{g}_t(x_t, l_t) \right]. \quad (34)$$

By the construction given in (30), we re-write the value function for the sequence of control laws $\hat{g}_{t+1:n-1}$ as

$$V_t(x_t, l_t) \leq c(x_t, \hat{g}_t(x_t, l_t)) + \mathbb{E} \left[V_{t+1}(X_{t+1}, \lambda_{t+1}^*) \mid x_t, \hat{g}_t(x_t, l_t) \right], \quad (35)$$

where, using (29), we recall that at time t the control law \hat{g}_t picks the optimal action u_t^* . Hence,

$$V_t(x_t, l_t) \leq c(x_t, u_t^*) + \mathbb{E} \left[V_{t+1}(X_{t+1}, \lambda_{t+1}^*) \mid x_t, \hat{g}_t(x_t, l_t) \right]. \quad (36)$$

Now, we can clearly see that the optimal value function at time t satisfies:

$$V_t(x_t, l_t) \leq \min_{\substack{u_t \in U_t \\ \lambda_{t+1} \in F_t(x_t, u_t, l_t)}} \left\{ c(x_t, u_t) + \mathbb{E} \left[V_{t+1}(X_{t+1}, \lambda_{t+1}(X_{t+1})) \mid x_t, u_t \right] \right\}. \quad (37)$$

The inequalities in (27) and (37) prove that the optimal value function at time t is both upper and lower bounded by the right-hand side of (13).

References

- Altman, E. (2021). *Constrained Markov decision processes*. Routledge.
- Ansere, J. A., Han, G., Liu, L., Peng, Y., & Kamal, M. (2020). Optimal resource allocation in energy-efficient internet-of-things networks with imperfect CSI. *IEEE Internet of Things Journal*, 7(6), 5401–5411.
- Baezner, M., & Robin, P. (2017). *Stuxnet: Tech. rep.*, ETH Zurich.

- Bertsekas, D., & Rhodes, I. (1973). Sufficiently informative functions and the minimax feedback control of uncertain dynamic systems. *IEEE Transactions on Automatic Control*, 18(2), 117–124.
- Chen, R. C., & Blankenship, G. L. (2004). Dynamic programming equations for discounted constrained stochastic control. *IEEE Transactions on Automatic Control*, 49(5), 699–709.
- Clement, J. G., & Kroer, C. (2021). First-order methods for wasserstein distributionally robust MDP. In *International conference on machine learning* (pp. 2010–2019). PMLR.
- Coraluppi, S. P., & Marcus, S. I. (2000). Mixed risk-neutral/minimax control of discrete-time, finite-state Markov decision processes. *IEEE Transactions on Automatic Control*, 45(3), 528–532.
- Dave, A., Chremos, I. V., & Malikopoulos, A. A. (2022). Social media and misleading information in a democracy: A mechanism design approach. *IEEE Transactions on Automatic Control*, 67(5), 2633–2639.
- Dave, A., Venkatesh, N., & Malikopoulos, A. A. (2022a). Approximate information states for worst-case control of uncertain systems. In *Proceedings of the 61th IEEE conference on decision and control* (pp. 4945–4950).
- Dave, A., Venkatesh, N., & Malikopoulos, A. A. (2022b). Decentralized control of two agents with nested accessible information. In *2022 American control conference* (pp. 3423–3430). IEEE.
- Dave, A., Venkatesh, N., & Malikopoulos, A. A. (2022c). On decentralized minimax control with nested subsystems. In *2022 American control conference* (pp. 3437–3444). IEEE.
- Dave, A., Venkatesh, N., & Malikopoulos, A. A. (2023). Approximate Information States for Worst-Case Control and Learning in Uncertain Systems. *arXiv:2301.05089* (in review).
- Ermon, S., Gomes, C., Selman, B., & Vladimirsky, A. (2012). Probabilistic planning with non-linear utility functions and worst-case guarantees. In *Proceedings of the 11th international conference on autonomous agents and multiagent systems-volume 2* (pp. 965–972).
- Gagrani, M., & Nayyar, A. (2017). Decentralized minimax control problems with partial history sharing. In *2017 American control conference* (pp. 3373–3379). IEEE.
- Ghiasi, M., Niknam, T., Wang, Z., Mehrandezh, M., Dehghani, M., & Ghadimi, N. (2023). A comprehensive review of cyber-attacks and defense mechanisms for improving security in smart grid energy systems: Past, present and future. *Electric Power Systems Research*, 215, Article 108975.
- Iyengar, G. N. (2005). Robust dynamic programming. *Mathematics of Operations Research*, 30(2), 257–280.
- Kim, K.-D., & Kumar, P. R. (2012). Cyber-physical systems: A perspective at the centennial. *Proceedings of the IEEE*, 100(Special Centennial Issue), 1287–1308.
- Malikopoulos, A. A. (2023a). On team decision problems with nonclassical information structures. *IEEE Transactions on Automatic Control*, 68(7), 3915–3930.
- Malikopoulos, A. A. (2023b). Separation of learning and control for cyber-physical systems. *Automatica*, 151(110912).
- Malikopoulos, A. A., Beaver, L. E., & Chremos, I. V. (2021). Optimal time trajectory and coordination for connected and automated vehicles. *Automatica*, 125(109469).
- Mannor, S., Simester, D., Sun, P., & Tsitsiklis, J. N. (2007). Bias and variance approximation in value function estimates. *Management Science*, 53(2), 308–322.
- Rasouli, M., Miehl, E., & Tenenetzis, D. (2018). A scalable decomposition method for the dynamic defense of cyber networks. In *Game theory for security and risk management* (pp. 75–98). Springer.
- Serror, M., Hack, S., Henze, M., Schuba, M., & Wehrle, K. (2020). Challenges and opportunities in securing the industrial internet of things. *IEEE Transactions on Industrial Informatics*, 17(5), 2985–2996.
- Shoukry, Y., Araujo, J., Tabuada, P., Srivastava, M., & Johansson, K. H. (2013). Minimax control for cyber-physical systems under network packet scheduling attacks. In *Proceedings of the 2nd ACM international conference on high confidence networked systems* (pp. 93–100).
- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT Press.
- Varapriya, V., & Singh, V. V. (2023). Chance-constrained formulation of MDPs under total reward criteria: an application to advertisement model. In *2023 European control conference* (pp. 1–6). IEEE.
- Venkatesh, N., Le, V.-A., Dave, A., & Malikopoulos, A. A. (2023). Connected and automated vehicles in mixed-traffic: Learning human driver behavior for effective on-ramp merging. In *Proceedings of the 62nd IEEE conference on decision and control* (pp. 92–97). IEEE.
- Wiesemann, W., Kuhn, D., & Rustem, B. (2013). Robust Markov decision processes. *Mathematics of Operations Research*, 38(1), 153–183.
- Zhang, L., Sridhar, K., Liu, M., Lu, P., Chen, X., Kong, F., et al. (2023). Real-time data-predictive attack-recovery for complex cyber-physical systems. In *2023 IEEE 29th real-time and embedded technology and applications symposium* (pp. 209–222).
- Zhu, Q., & Başar, T. (2011). Robust and resilient control design for cyber-physical systems with an application to power systems. In *2011 50th IEEE conference on decision and control and European control conference* (pp. 4066–4071). IEEE.