# Adversarial Examples for Model-Based Control: A Sensitivity Analysis

Po-han Li[1], Ufuk Topcu[2], and Sandeep P. Chinchali[1]

*Abstract—*

We propose a method to attack controllers that rely on external timeseries forecasts as task parameters. An adversary can manipulate the costs, states, and actions of the controllers by forging the timeseries, in this case perturbing the real timeseries. Since the controllers often encode safety requirements or energy limits in their costs and constraints, we refer to such manipulation as an adversarial attack. We show that different attacks on model-based controllers can increase control costs, activate constraints, or even make the control optimization problem infeasible. We use the linear quadratic regulator and convex model predictive controllers as examples of how adversarial attacks succeed and demonstrate the impact of adversarial attacks on a battery storage control task for power grid operators. As a result, our method increases control cost by $8500\%$ and energy constraints by $13\%$ on real electricity demand timeseries.

## I. Introduction

There are rich applications for model predictive controllers (MPC) that rely on timeseries forecasts as task parameters. For example, cellular network traffic schedulers predict city-wide mobility data to assign base station connections to mobile devices [1], power grid operators use electricity demand patterns to optimize battery storage [2], [3], and stock traders use price forecasts to make trading decisions. In these applications, the timeseries are not measured nor determined by the controllers, but by external sources. We refer to such controllers as ***input-driven controllers***, where controllers use reliable estimates of internal control states and dynamics as well as external timeseries forecasts to make decisions, known as actions or controls. Since the controller plays a passive role in receiving external timeseries, a natural question is: are the timeseries forecasts also reliable? In this paper, we use the linear quadratic regulator (LQR) to discuss how epistemic uncertainty or malicious external sources can affect control cost or constraints. We further extend the discussion to convex MPC controllers.

**Related work:** Our system model is close to [2], which proposes an input-driven LQR controller with external forecasts of timeseries. However, in contrast to our work, it focuses on the optimal compression for timeseries across a bandwidth-limited network, while we instead focus on adversarial attacks.

Adversarial attacks make bounded, often human-imperceptible, perturbations on a sensory input (e.g., image) to cause errors in output predictions (i.e., image classifications). [4] studies adversarial attacks

[1]Department of Electrical and Computer Engineering, The University of Texas at Austin, [2]Oden Institute for Computational Engineering and Sciences, The University of Texas at Austin, {pohanli, utopcu, sandeepc}@utexas.edu
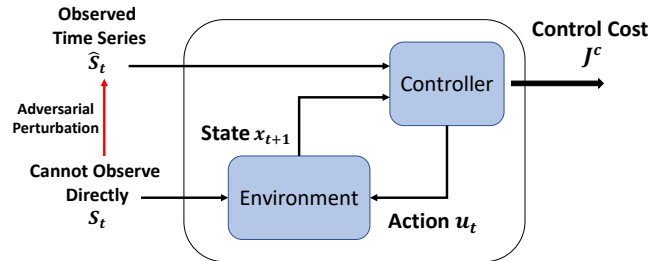
Fig. 1: **Adversarial Attacks on Timeseries Forecasts For Model-Based Control.** Many modern controllers require reliable forecasts of demand or prices to make decisions. In this paper, we show how slight perturbations in a forecast can dramatically increase control costs or violate control constraints. These errors in forecasting can occur due to out-of-distribution (OoD) timeseries, natural noise, or adversarial perturbations. Specifically, at time $t$, a controller observes state $x_t$ and timeseries $\hat{S}_t$, which is perturbed by an adversarial source. Then, it takes action $u_t$ to minimize the cost $J^c$. The next state $x_{t+1}$ is determined by the real timeseries $S_t$, action $u_t$, and previous state $x_t$. The adversarial source perturbs timeseries $S_t$ to $\hat{S}_t$ within a bounded perturbation in order to increase the control cost or make the controller violate constraints.

on probabilistic autoregressive models, and [5], [6] focus on adversarial noise for image classification. While [7], [8], [9], [10], [11] study adversarial attacks that affect the dynamics of a reinforcement learning (RL) agent, our work exploits the structure of a model-based control task to generate adversarial attacks on timeseries inputs.

Adversarial attacks in control systems have also been studied in [12], [13], [14]. They formulate the adversarial attack of a controller as a max-min problem, where the adversary's goal is to maximize the cost and the controller's goal is to minimize it. In [12], the adversary adds perturbations to the data set of a non-linear data-driven controller, while perturbations are a type of noise which can affect the controller's states directly in [13], [14]. Our work is distinct since we focus on how an adversary can perturb external forecasts of timeseries to (a) maximize the control cost and (b) make a controller nearly violate its strict state and control constraints. To the best of our knowledge, ours is one of the first works to describe how to introduce attacks that make a controller nearly violate state and control constraints, which are important in practice, for example, to express energy limits.

**Insights and Contributions:** Our key technical insight is that input-driven controllers are sensitive to errors in external timeseries, and input-driven controllers have different angles of attack, i.e. increasing control costs, activating constraints, or even making the control optimization problem infeasible. Such attacks are important since constraints are often es-

sential to control tasks – state constraints describe desired safety levels or operating regions, and action constraints can be power or energy constraints. Therefore, we address possible white-box attacks, where a malicious external source knows the controller's parameters and its dynamics. Also, we formulate all attacks as a bounded perturbation on real timeseries, since the adversarial timeseries should be similar to the original one in order to be indistinguishable by human-guided analysis or anomaly detectors.

Based on these insights, our contributions are three-fold. First, we analytically calculate the optimal, bounded perturbation of a timeseries forecast in order to increase a model-based controller's cost. Second, we provide a numerical method to generate adversarial timeseries that attack strict control constraints. Lastly, we show by numerical experiments on real electricity demand data that adversarial attacks for control costs differ from those that violate strict state or control constraints.

## II. System Model and Problem Formulation

### A. System model

Imagine a battery storage operator's controller must decide whether to charge or discharge its batteries based on electricity price forecasts in the market. At time $t$, the electricity demand is an external timeseries $S_t \in \mathbb{R}^p$, which cannot be affected by the controller. The state $x_t \in \mathbb{R}^n$ represents the charge on batteries and action $u_t \in \mathbb{R}^m$ represents how much to charge the batteries in order to minimize the cost. The controller cannot affect the external timeseries $S_t$ since it is independent of actions $u_t$ and forecasted by an external source. We denote **full future** control vectors in bold fonts specifically, $\boldsymbol{u} = u_{0:T-1} \in \mathbb{R}^{mT}$, $\boldsymbol{S} = S_{0:T-1} \in \mathbb{R}^{pT}$, and $\boldsymbol{x} = x_{0:T} \in \mathbb{R}^{m(T+1)}$ for a finite time horizon $T$. The system dynamics is therefore determined by the external timeseries $\boldsymbol{S}$, action $\boldsymbol{u}$, and state $\boldsymbol{x}$. The controller may not be able to observe the real timeseries $\boldsymbol{S}$. Due to measurement noise, forecasting error, or attacks from an external source, it can only observe a perturbed timeseries $\hat{\boldsymbol{S}}$. However, we assume the controller perfectly knows its internal plant state $x_t$.

The system model is shown in Fig. 1. First, at time $t$, the controller observes the perturbed timeseries $\hat{S}_t$ and current state $x_t$ and then decides an action $u_t$. Second, the action $u_t$, current state $x_t$, and real timeseries $S_t$ determine the next state $x_{t+1}$ through the controller's known dynamics. The controller aims to minimize its control cost $J^c(\boldsymbol{u}; \boldsymbol{S}, x_0)$, which is a function of all actions $\boldsymbol{u}$, initial state $x_0$, and real timeseries $\boldsymbol{S}$. Note that the optimal action of a controller $\boldsymbol{u}^*(x_0; \hat{\boldsymbol{S}})$ is a function of initial state $x_0$, and the observed timeseries $\hat{\boldsymbol{S}}$. We next formulate the problem for a special case of the linear quadratic regulator (LQR).

### B. Problem Formulation

We now describe the system dynamics of input-driven LQR and derive its control cost to obtain the sensitivity of the control cost with respect to forecasting errors. Our derivation extends the work of [2], which introduces input-driven LQR but does not address adversarial attacks. For

clarity, we summarize the derivation here and refer readers to [2] for details. For any time step $t$, the linear system dynamics are given by:

$$x_{t+1} = \boldsymbol{A}x_t + \boldsymbol{B}u_t + \boldsymbol{C}S_t,$$

where $\boldsymbol{A} \in \mathbb{R}^{n \times n}$, $\boldsymbol{B} \in \mathbb{R}^{n \times m}$, and $\boldsymbol{C} \in \mathbb{R}^{n \times p}$ are the parameters describing how previous state $x_t$, action $u_t$, and external timeseries $S_t$ affect the next state $x_{t+1}$. We rewrite the system dynamics in a non-recursive form:

$$x_{t+1} = \boldsymbol{A}^{t+1}x_0 + \boldsymbol{M}_t\boldsymbol{u} + \boldsymbol{N}_t\boldsymbol{S}, \tag{1}$$

where $\boldsymbol{M}_t = [\boldsymbol{A}^t\boldsymbol{B} \ \ \boldsymbol{A}^{t-1}\boldsymbol{B} \ \ ... \ \ \boldsymbol{B} \ \ \boldsymbol{0}] \in \mathbb{R}^{n \times mT}$, $\boldsymbol{N}_t = [\boldsymbol{A}^t\boldsymbol{C} \ \ \boldsymbol{A}^{t-1}\boldsymbol{C} \ \ ... \ \ \boldsymbol{C} \ \ \boldsymbol{0}] \in \mathbb{R}^{n \times pT}$. The linear quadratic cost function is defined as:

$$J^c(\boldsymbol{u}; \boldsymbol{S}, x_0) = \sum_{t=0}^{T} x_t^\top \boldsymbol{Q}x_t + \sum_{t=0}^{T-1} u_t^\top \boldsymbol{R}u_t, \quad \boldsymbol{Q} \succ 0, \boldsymbol{R} \succ 0. \tag{2}$$

Positive-definite matrices $\boldsymbol{Q}$ and $\boldsymbol{R}$ are represented as $\boldsymbol{Q}, \boldsymbol{R} \succ 0$. The control cost is a function of $\boldsymbol{S}$ and $\boldsymbol{u}$

$$J^c(\boldsymbol{u}; \boldsymbol{S}, x_0) = \boldsymbol{u}^\top \underbrace{\left( \text{BlockDiag}(\boldsymbol{R}, T) + \sum_{t=0}^{T-1} \boldsymbol{M}_t^\top \boldsymbol{Q}\boldsymbol{M}_t \right)}_{\boldsymbol{K}} \boldsymbol{u}$$

$$+ 2 \underbrace{\left[ \sum_{t=0}^{T-1} \boldsymbol{M}_t^\top \boldsymbol{Q}(\boldsymbol{A}^{t+1}x_0 + \boldsymbol{N}_t\boldsymbol{S}) \right]^\top}_{\boldsymbol{k}(x_0, s)^\top} \boldsymbol{u} + \text{constant term},$$
$$\tag{3}$$

where $\text{BlockDiag}(\boldsymbol{R}, T) \in \mathbb{R}^{mT \times mT}$ is a block matrix placing $T$ $\boldsymbol{R}$ matrices on the diagonal, and the constant term independent of $\boldsymbol{u}$ is $\sum_{t=0}^{T-1} \left(\boldsymbol{A}^{t+1}x_0 + \boldsymbol{N}_t\boldsymbol{S}\right)^\top \boldsymbol{Q} \left(\boldsymbol{A}^{t+1}x_0 + \boldsymbol{N}_t\boldsymbol{S}\right)$. By Eq. 3, optimal actions are determined by the external timeseries which the controller observes:

$$\boldsymbol{u}^*(x_0; S) = \arg\min_{\boldsymbol{u}} J^c(\boldsymbol{u}; \boldsymbol{S}, x_0) = -\boldsymbol{K}^{-1}k(x_0, \boldsymbol{S}).$$

Thus, we obtain different optimal actions based on different observed timeseries and show their difference:

$$\boldsymbol{u}^* = -\boldsymbol{K}^{-1}k(x_0, \boldsymbol{S}), \ \hat{\boldsymbol{u}}^* = -\boldsymbol{K}^{-1}k(x_0, \hat{\boldsymbol{S}}). \tag{4}$$

$$\hat{\boldsymbol{u}}^* - \boldsymbol{u}^* = -\boldsymbol{K}^{-1} \left( k(x_0, \hat{\boldsymbol{S}}) - k(x_0, \boldsymbol{S}) \right)$$

$$= -\boldsymbol{K}^{-1} \underbrace{\sum_{t=0}^{T-1} \boldsymbol{M}_t^\top \boldsymbol{Q}\boldsymbol{N}_t}_{\boldsymbol{L}}(\hat{\boldsymbol{S}} - \boldsymbol{S}). \tag{5}$$

Eq. 5 shows the errors in control are linear with respect to errors in forecasting, and the coefficient is determined by the parameters of the system dynamics and cost function. Next, we focus on the sensitivity of control cost $\Delta J$, defined as follows, and show that it is quadratic in forecasting error. That is, different elements of timeseries in different time

steps have nonidentical effects on the change in control cost:

$$\Delta J^c = J^c(\hat{\boldsymbol{u}}^*; \boldsymbol{S}, x_0) - J^c(\boldsymbol{u}^*; \boldsymbol{S}, x_0)$$
$$= (\hat{\boldsymbol{u}}^* - \boldsymbol{u}^*)^\top \boldsymbol{K} (\hat{\boldsymbol{u}}^* - \boldsymbol{u}^*) \qquad (6)$$
$$= (\hat{\boldsymbol{S}} - \boldsymbol{S})^\top \underbrace{\boldsymbol{L}^\top \boldsymbol{K}^{-1} \boldsymbol{L}}_{\boldsymbol{\Psi}} (\hat{\boldsymbol{S}} - \boldsymbol{S}).$$

Note that $\boldsymbol{\Psi} \succeq 0$ since $\boldsymbol{K} \succ 0$. A way to intuitively explain Eq. 6 is that the change in control cost $\Delta J^c$ is determined by the matrix $\boldsymbol{\Psi}$, which weights the errors in elements of $S$ based on their importance to the control cost. In the next section, we discuss different approaches to attack the controller by perturbing the timeseries $S$ to $\hat{S}$.

## III. ADVERSARIAL ATTACKS ON CONTROL COST AND CONSTRAINTS

### A. Cost Adversarial Attack

We now derive the optimal attack $\hat{\boldsymbol{S}}$ for an input-driven LQR controller given real timeseries $\boldsymbol{S}$ and a perturbation bound $\delta$. We define $\delta$ as the upper bound of the perturbation in the $L_2$ norm. The perturbation is restricted within a bound since attackers aim to generate timeseries similar to the original that are indistinguishable by human inspection or automated anomaly detectors. The problem is formulated as:

$$\max_{\hat{\boldsymbol{S}}} \quad \Delta J^c = (\hat{\boldsymbol{S}} - \boldsymbol{S})^\top \boldsymbol{\Psi} (\hat{\boldsymbol{S}} - \boldsymbol{S}).$$
$$\text{subject to} \quad \|\hat{\boldsymbol{S}} - \boldsymbol{S}\|_2 \le \delta \qquad (7)$$

*Theorem 1 (Optimal Perturbation of Cost):* The optimal timeseries maximizing $\Delta J^c$ is:

$$\hat{\boldsymbol{S}}^* = \arg\max_{\hat{\boldsymbol{S}}} \quad (\hat{\boldsymbol{S}} - \boldsymbol{S})^\top \boldsymbol{\Psi} (\hat{\boldsymbol{S}} - \boldsymbol{S}) = \boldsymbol{S} \pm \delta v_1.$$
$$\text{subject to} \quad \|\hat{\boldsymbol{S}} - \boldsymbol{S}\|_2 \le \delta \qquad (8)$$

The corresponding change in control cost is:

$$\max_{\hat{\boldsymbol{S}}} \Delta J^c = \delta^2 \lambda_1, \qquad (9)$$

where $\lambda_1$, $v_1$ are the dominant (largest) eigenvalue and eigenvector of $\boldsymbol{\Psi}$.

*Proof:* We rewrite the constraint in Eq. 7 as:

$$\|\hat{\boldsymbol{S}} - \boldsymbol{S}\|_2^2 \le \delta^2.$$

Then, using the method of Lagrange Multipliers:

$$\nabla_{(\hat{\boldsymbol{S}} - \boldsymbol{S})} \Delta J^c = \lambda \nabla_{(\hat{\boldsymbol{S}} - \boldsymbol{S})} (\delta^2 - \|\hat{\boldsymbol{S}} - \boldsymbol{S}\|_2^2).$$

Using the derivative of a quadratic functions, we see

$$\boldsymbol{\Psi}(\hat{\boldsymbol{S}} - \boldsymbol{S}) = -\lambda(\hat{\boldsymbol{S}} - \boldsymbol{S}).$$

So far, we showed the optimal solution $\hat{\boldsymbol{S}}^*$ occurs when $\hat{\boldsymbol{S}} - \boldsymbol{S}$ is $\boldsymbol{\Psi}$'s eigenvector. Since we aim to maximize $\Delta J^c$, $\hat{\boldsymbol{S}} - \boldsymbol{S}$ is chosen as the dominant eigenvector $v_1$ of $\boldsymbol{\Psi}$ corresponding to the largest eigenvalue $\lambda_1$. Without loss of generality, we set $\|v_1\|_2 = 1$.

Lastly, the corresponding optimal solutions and value are:

$$\hat{\boldsymbol{S}} = \boldsymbol{S} \pm \delta v_1, \quad \max_{\hat{\boldsymbol{S}}} \Delta J^c = \delta^2 \lambda_1.$$

The result is two-fold. First, there are two optimal solutions and two optimal perturbations. Second, the optimal perturbation is independent of the real timeseries and the initial state, but only depends on the system dynamics parameters and the cost function. That is, the adversary only needs to know the system dynamics parameters and cost metrics of the controller, not the real timeseries or the initial state of the controller, thus making the attack easier to implement.

### B. Control-agnostic Attack is Random

We consider a benchmark where the attack has the same bound on perturbations, but is agnostic to the control task. Namely, the perturbation is determined without any knowledge of the system dynamics and the cost function. Similar to Eq. 7, the problem is formulated as:

$$\max_{\hat{\boldsymbol{S}}} \quad \Delta J^c = (\hat{\boldsymbol{S}} - \boldsymbol{S})^\top (\hat{\boldsymbol{S}} - \boldsymbol{S}).$$
$$\text{subject to} \quad \|\hat{\boldsymbol{S}} - \boldsymbol{S}\|_2 \le \delta \qquad (10)$$

This is a special case of Thm. 1 where $\boldsymbol{\Psi}$ is the identity matrix $\mathbf{I}$. In this case, all unit vectors are eigenvectors, and the corresponding eigenvalue is always 1. Consequently, all vectors which activate constraint $\|\hat{\boldsymbol{S}} - \boldsymbol{S}\|_2 = \delta$ are optimal, and the optimal value is $\delta^2$.

### C. Constraint Adversarial Attack and A More General Form

In the previous sections, we only consider the simplest case of an input-driven LQR controller without constraints. We now consider a more general setting of an input-driven controller with linear dynamics, which is a ***convex optimization problem***:

$$\min_{\boldsymbol{u}} \quad J^c(\boldsymbol{u}; x_0, \boldsymbol{S}) \qquad (11a)$$
$$\text{subject to} \quad x_{t+1} = \boldsymbol{A}^{t+1} x_0 + \boldsymbol{M}_t \boldsymbol{u} + \boldsymbol{N}_t \boldsymbol{S},$$
$$\forall t = 0, \dots, T-1 \qquad (11b)$$
$$f_i(\boldsymbol{u}; x_0, \boldsymbol{S}) \le 0, \; \forall i = 1, \dots, n_{\text{ineq}} \qquad (11c)$$
$$g_i(\boldsymbol{u}; x_0, \boldsymbol{S}) = 0, \; \forall i = 1, \dots, n_{\text{eq}}, \qquad (11d)$$

where 11a is a convex cost function, and 11b is the system dynamics. 11c and 11d are constraints of the controller. The optimal solution of Eq. 11 is defined as $\boldsymbol{u}^*_{\text{gen}}$. In most of the cases, $\boldsymbol{u}^*_{\text{gen}}$ does not have an analytical solution and one can only obtain a numerical solution using convex optimization solvers such as [15], [16], [17], [18]. Note that all $f(\boldsymbol{u}; x_0, \boldsymbol{S})$ are convex and $g(\boldsymbol{u}; x_0, \boldsymbol{S})$ are affine since this is a convex optimization problem. $\boldsymbol{u}^*_{\text{gen}}$ is determined by parameters $x_0$ and $\boldsymbol{S}$ given fixed $f$ and $g$. When the optimization problem is defined with a set of parameters $x_0$ and $\boldsymbol{S}$, the solution map is defined as a function mapping the parameters to the solution. We thus rewrite $\boldsymbol{u}^*_{\text{gen}}$ as a solution map $\boldsymbol{u}^*_{\text{gen}}(x_0; \boldsymbol{S})$.

We now formulate the optimal attack to the general setting of Eq. 11 then approximate the attack by a first-order Taylor Series expansion. As before, the controller only observes the

perturbed timeseries $\hat{\boldsymbol{S}}$, and the perturbation aims to maximize any differentiable function $h_{\mathrm{adv}}(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); x_0, \boldsymbol{S})$ within an upper bound $\delta$. $h_{\mathrm{adv}}(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); x_0, \boldsymbol{S})$ is referred to as the target function. Then, the optimal attack is:

$$\hat{\boldsymbol{S}}^* = \arg\max_{\hat{\boldsymbol{S}}} \quad h_{\mathrm{adv}}(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); x_0, \boldsymbol{S}). \quad (12a)$$

$$\text{subject to} \quad \|\hat{\boldsymbol{S}} - \boldsymbol{S}\|_2 \leq \delta \quad (12b)$$

When the target function $h_{\mathrm{adv}}$ is the change in control cost $\Delta J^c$ as in Eq. 6, and the solution map of actions are the same as Eq. 4, namely,

$$
\begin{aligned}
h_{\mathrm{adv}}(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); x_0, \boldsymbol{S}) &= \Delta J^c \\
&= J^c(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); \boldsymbol{S}, x_0) \\
&\quad - J^c(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \boldsymbol{S}); \boldsymbol{S}, x_0) \\
\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}) &= -\boldsymbol{K}^{-1}k(x_0, \hat{\boldsymbol{S}}) \\
\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \boldsymbol{S}) &= -\boldsymbol{K}^{-1}k(x_0, \boldsymbol{S}),
\end{aligned}
\quad (13)
$$

the problem is identical to the one formulated in Thm. 1.

Eq. 12 may not be a convex optimization problem, since solution maps of a convex optimization problem are not concave nor convex in general. However, one can approximately solve Eq. 12 by $\nabla_{\hat{\boldsymbol{S}}}\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}})$, the gradient of $\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}})$ with respect to $\hat{\boldsymbol{S}}$, as described in [19], [20]. Here, $\nabla_{\hat{\boldsymbol{S}}}\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}) \in \mathbb{R}^{pT \times mT}$ is a matrix where its $(i,j)$ element is $\partial u_{\mathrm{gen},j}^*/\partial \hat{s}_i$ with subscripts $i, j$ denoting the $i, j$ elements of $\hat{\boldsymbol{S}}_i$ and $\boldsymbol{u}_{\mathbf{gen},j}^*$, respectively. Hence, we can solve Eq. 12 using $\nabla_{\hat{\boldsymbol{S}}}\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}})$ as a linear local approximation. Assuming that $\delta$ is small, and $h_{\mathrm{adv}}$ is differentiable with respect to $\boldsymbol{u}_{\mathbf{gen}}^*$, Eq. 12 can be approximated as:

$$
\begin{aligned}
\hat{\boldsymbol{S}}^* &\approx \boldsymbol{S} + \delta \times \mathrm{Unit}(\nabla_{\hat{\boldsymbol{S}}}h_{\mathrm{adv}}(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); x_0, \boldsymbol{S})|_{\hat{\boldsymbol{S}}=\boldsymbol{S}}) \\
&= \boldsymbol{S} + \delta \times \mathrm{Unit}(\nabla_{\hat{\boldsymbol{S}}}\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}})|_{\hat{\boldsymbol{S}}=\boldsymbol{S}} \times \\
&\quad \begin{bmatrix} \partial h_{\mathrm{adv}}/\partial u_{\mathrm{gen},1}^* \\ \partial h_{\mathrm{adv}}/\partial u_{\mathrm{gen},2}^* \\ \vdots \\ \partial h_{\mathrm{adv}}/\partial u_{\mathrm{gen},mT}^* \end{bmatrix} |_{\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \boldsymbol{S})}) \text{(chain rule)}
\end{aligned}
$$
$$(14)$$

Here, $\mathrm{Unit}(\cdot)$ is the operator normalizing a vector to a unit vector.

**Constraint Adversarial $h_{\mathrm{adv}}$:** Eq. 12 and 14 show how attackers can maximize any differentiable function, but we focus on cases where the target function is related to the inequality constraints in Eq. 11. When the attack target is related to the inequality constraints, the attacker can bring the controller closer to activating or even violating its constraints. Taking the battery storage operator's controller as an example, an attacker can increase the controller's energy consumption when the controller has energy constraints. Also, the cost function may not capture the constraints, thus making it more difficult to detect the presence of attacks. Note that in some cases, when all constraints of a MPC problem are violated, the control problem may become infeasible. We list some common scenarios and the corresponding

constraints here.

1) ***Box-constrained LQR*** There are upper and lower constraints of control actions, $u_{min} \leq \boldsymbol{u} \leq u_{max}$. Then

$$h_{\mathrm{adv}}(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); x_0, \boldsymbol{S}) = \max_t\{u_{\mathrm{gen,t}}^*\}, \text{ or}$$

$$h_{\mathrm{adv}}(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); x_0, \boldsymbol{S}) = \max_t\{-u_{\mathrm{gen,t}}^*\}.$$

2) ***Energy Limitation*** The controller has an energy constraint in the finite time horizon $T$, so its $L_1$ norm of actions is limited, Namely, $\|\boldsymbol{u}_{\mathbf{gen}}^*\|_1 - \text{constant} \leq 0$. Then,

$$h_{\mathrm{adv}}(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); x_0, \boldsymbol{S}) = \|\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}})\|_1.$$

**Extensions and Limitations:** Eq. 14 uses a gradient to calculate the first order Taylor expansion of the solution map. One can also use more sophisticated gradient ascent methods, such as the optimizer ADAM [21], to update $\hat{\boldsymbol{S}}$ as long as Eq. 12b holds. Also Eq. 12b need not be the $L_2$ norm, so Eq. 14 can be modified to other norms by changing the $\mathrm{Unit}(\cdot)$ operator to normalize other norms. However, the approximation addressed in Eq. 14 has a shortcoming – when any gradient is $\boldsymbol{0}$, $\hat{\boldsymbol{S}}^* = \boldsymbol{S}$. Approximation fails in this case and returns the same timeseries. It happens when $h_{\mathrm{adv}}(\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}}); x_0, \boldsymbol{S})$ or $\boldsymbol{u}_{\mathbf{gen}}^*(x_0; \hat{\boldsymbol{S}})$ have local extrema, and Eq. 8 is exactly the case. The reason is simple– we are only using first order approximation to solve Eq. 12b. If one can calculate higher order terms, this shortcoming will be resolved. Nevertheless, second order differentiation of a convex optimization problem is an open research topic and out of the scope of this paper.

## IV. EXPERIMENTS

The goal of our experiments is to show that we can slightly perturb a timeseries to targetedly attack control cost and control constraints. We show that our targeted adversarial attacks affect control cost significantly more than random perturbations of the timeseries with the same perturbation bound, illustrating the efficacy and potency of the attacks.

We evaluate our methods on two sets of timeseries. The first one is a synthetic Autoregressive Integrated Moving Average (ARIMA) process [22] generated by random parameters. The other is hourly electric energy consumption data from PJM Interconnection LLC [23], a regional transmission organization (RTO) in the Eastern United States. We sample the timeseries from the data set of American Electric Power from 2015. In each set, a distribution of timeseries is input to the fixed control task described above. For both sets, $n = m = p = 1$, $x_0 = 1$, $\boldsymbol{A} = \boldsymbol{C} = \boldsymbol{Q} = \boldsymbol{R} = 1$, and $\boldsymbol{B} = -1$. $T = 50$ for ARIMA, and $T = 120$ for PJM electricity. ARIMA's experiment is purely synthetic, while the PJM electricity tasks is to simulate real world battery operation. The operator decides how much to discharge or charge the battery (action $u_t$) so that the stored battery capacity (state $x_t$) can meet the electricity market demand forecast (external timeseries). The operator's goal is to minimize the control
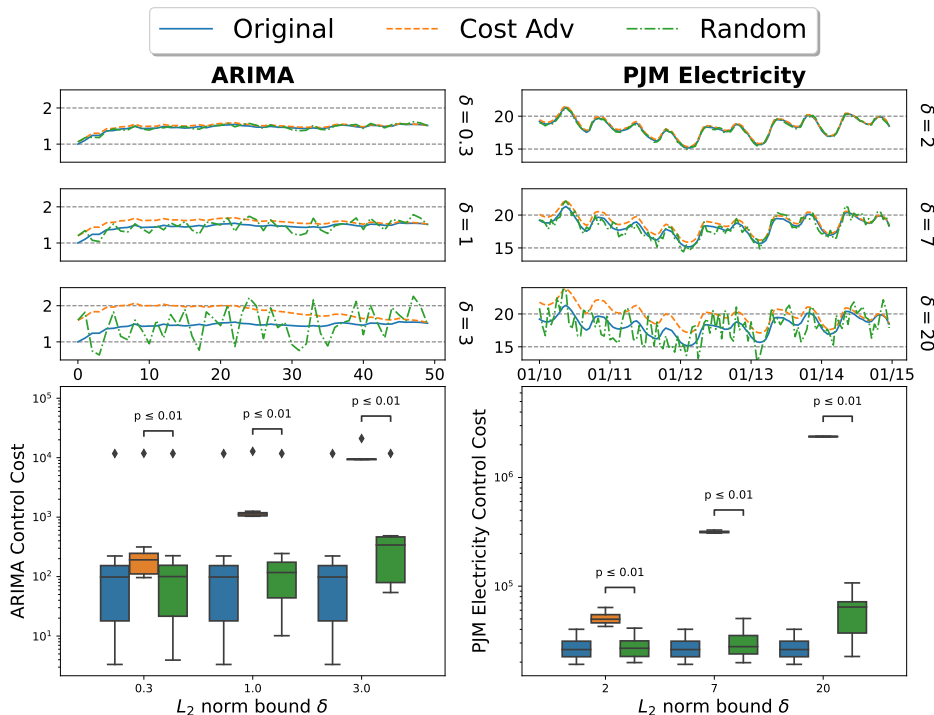
Fig. 2: **Adversarial Timeseries can affect control cost dramatically.** The first three rows show how the original and perturbed timeseries differ under various perturbation bounds $\delta$. The last row shows the corresponding control costs in different scenarios on a log scale. In general, the dominant eigenvalue $\lambda_1$ is large for high dimensional $\mathbf{\Psi}$, and control costs of *Cost Adv* are quadratic in the perturbation bound $\delta$ and linear to the eigenvalue $\lambda_1$ (Thm. 1). Consequently, although the timeseries look similar, the resulting costs are increased by orders of magnitude, illustrating that our synthesized attacks are indeed powerful. Indeed, the adversarial attacks significantly increase control cost compared to the *Original* timeseries and a benchmark of random perturbations with the same bound (*Random*) with a significant Wilcoxon p-value of $\leq 0.01$.

cost, expressed as a quadratic combination of the current states and actions. The operator wants the battery capacity to reach a set-point of half-full at every time step, represented as $x_t = 0$, to allow flexibly switching between favorable markets according to [3]. Therefore, the control cost matrix $\mathbf{Q}$ weights a penalty on the state's deviation over or under the set-point. Also, the cost matrix $\mathbf{R}$ penalizes the amount of charging and discharging to preserve battery health.

We first quantitatively show how much bounded adversarial perturbations $\hat{\mathbf{S}}$ can increase the control cost $J^c(\mathbf{u}^*_{gen}(x_0; \hat{\mathbf{S}}); \mathbf{S}, x_0)$ in Fig.2. Next, we show how perturbed timeseries can affect the control actions $\mathbf{u}$ with different target functions such as $\max\{u_t\}$ and $\|\mathbf{u}\|_1$ in Fig. 3. We compare our methods with a benchmark of adding random noise to the timeseries, denoted as *Random* in the figures. The x-axes of ARIMA and PJM electricity timeseries are time steps and date (in hours) in Fig. 2 and 3, respectively.

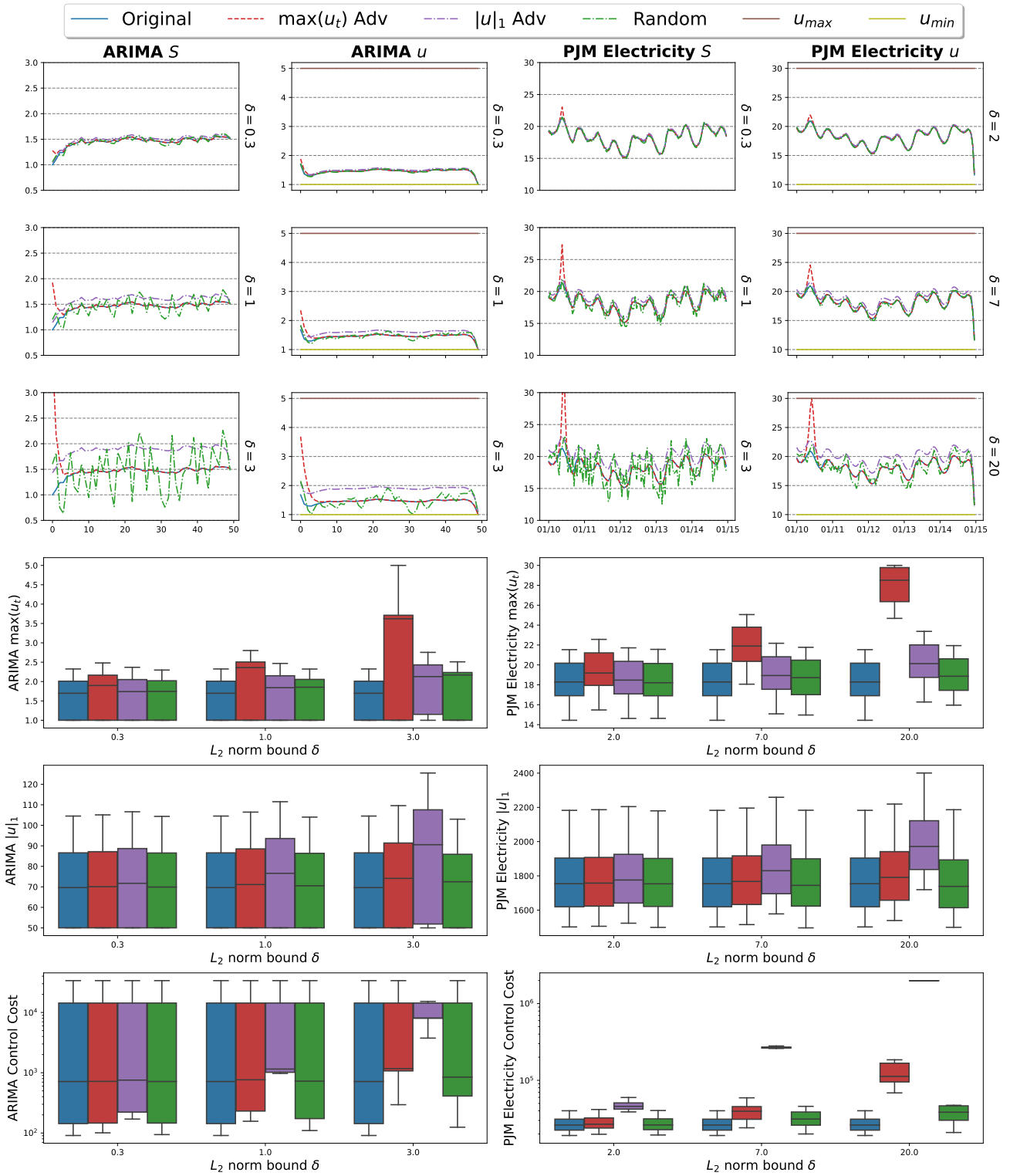### A. Adversarial Attacks on Control Cost

In Fig. 2, we compare the two sets of timeseries to the original one under two attack scenarios: Cost adversarial (derived in Thm. 1) and Random attack, denoted as *Cost Adv* and *Random*, respectively. Note that random attack is the same as the control-agnostic perturbation as described in Eq. 10 since all vectors are dominant eigenvectors for control-agnostic. As shown in Fig. 2, when the perturbation bound $\delta$ increases, *Cost Adv* and *Random* deviate more from the original timeseries. In general, the matrix determining the cost $\mathbf{\Psi}$ has a large dominant eigenvalue $\lambda_1$ when its dimension is large, so adversarial control costs $J^c(\mathbf{u}^*_{gen}(x_0; \hat{\mathbf{S}}); \mathbf{S}, x_0)$ dramatically increase. When the timeseries perturbation bounds $\delta$ are 0.3, 1 and 3, control costs increase 7.4%, 82%, and 740% on average for ARIMA. Similarly, when the perturbation bounds $\delta$ are 2, 7 and 20,

control costs increase 85%, 1000%, and 8500% on average for PJM electricity. The control cost is quadratic in the perturbation bound $\delta$, which is consistent with Thm. 1.

In the PJM electricity experiments, the external attack can affect the cost by orders of magnitude, which reflects the flexibility of the operator to switch markets as well as battery health. We also show the Wilcoxon p-values [24] of cost distributions between *Cost Adv* and *Random*. All p-values are statistically significant (less than 0.01) and thus show the efficacy of the *Cost Adv* method compared to random perturbations. Clearly, adversarial timeseries similar to original ones can increase the control costs dramatically by orders of magnitude for both ARIMA and PJM electricity datasets.

### B. Constraint Adversarial

In Fig. 3, we show how timeseries can affect the control actions and the resulting control cost $J^c(\mathbf{u}^*_{gen}(x_0; \hat{\mathbf{S}}); \mathbf{S}, x_0)$. The controller has constraints on the maximum and minimum values of the action, denoted as $u_{max}$ and $u_{min}$ in Fig. 3. We compare three attack scenarios with the original timeseries: 1. Box-constrained LQR. 2. Energy Limitation. 3. Random. The target functions of Box-constrained LQR and Energy Limitation are maximum action $\max\{u_t\}$ and absolute sum of actions $\|\mathbf{u}\|_1$, as described in Sec. III-C. In Fig. 3, when the perturbation bound $\delta$ increases, all three perturbed timeseries deviate more from the original one. The adversarial timeseries maximizing the maximum action $\max\{u_t\}$ only perturbs a small interval of time steps since its target function only relates to some consecutive elements of the control actions. On the other hand, the timeseries maximizing the absolute sum of actions $\|\mathbf{u}\|_1$ shifts every time step since its target function relates to the absolute

Fig. 3: **Adversarial attacks that target control cost and control constraints are inherently different.** The first three rows show how timeseries and resulting control actions differ under various perturbation bounds $\delta$. $\max(u_t)$ aims to maximize the maximum value of actions, so it only perturbs the timeseries in a short time interval where the maximum value originally occurs. On the other hand, $\|u\|_1$ aims to maximize the absolute sum of control actions so it perturbs almost every time step. The resulting maximum values of actions, absolute sums of actions, and control costs are shown in the last three rows of the plot. For the first two rows of box plots, when the perturbation bound $\delta$ increases, the corresponding attacks increase the target functions more while other scenarios are closer to the original one. For $\max(u_t)$ of PJM electricity , when the perturbation bound $\delta$ is 20, all cases reach the maximum constraint value of the action, so it is a horizontal line not a boxed distribution. The last row shows that control costs are not necessarily increased under different attacks that target control constraints, illustrating our main point that the two targets of attack are indeed different. Note that control costs are shown on a log scale, so $\max(u_t)$ of PJM electricity appears as a narrow distribution. Since attacks on different targets of control cost and control constraints are independent, we show that there exist diverse attack surfaces on controllers.

sum of all actions.

The fourth row of Fig. 3 shows the resulting maximum value of actions in different scenarios and perturbation bounds $\delta$. $\max\{u_t\}$ increases $4.6\%$, $15\%$, and $47\%$ on average for ARIMA and $5.5\%$, $20\%$, and $53\%$ on average for PJM electricity under perturbation bounds $0.3$, $1$, $3$, and $2$, $7$, $20$, respectively. As for $\|\boldsymbol{u}\|_1$ in the fifth row of Fig. 3, it increases $1.2\%$, $4\%$, and $13\%$ on average for ARIMA and $1.2\%$, $4.3\%$, and $12\%$ on average for PJM electricity under perturbation bounds $0.3$, $1$, $3$, and $2$, $7$, $20$, respectively. In the last row of Fig. 3, we show that different target functions may increase the control costs, but not as significantly as Fig. 2, since they are not directly related to the target function of the attack. Hence, we do not show the Wilcoxon p-values.

The results confirm that the attack for a specific target function is significantly different than for other target functions. In the PJM electricity experiments, external attacks can affect the maximum value or absolute sum of the operator's actions to discharge its battery. The attacks cause the controller to discharge more power to the battery, thus potentially reducing its operating life, or consuming more energy when charging the battery.

## V. Conclusion and Future Work

In this paper, our key insight is that input-driven controllers are sensitive to external attacks, and there are different attack targets besides control cost, such as control constraints. Thus, we propose a general formulation of white-box attacks, where attackers can perturb the external timeseries observed by the controller to maximize any differentiable function that may or may not be the control cost. Our key contribution is to first show an analytical form of an attack on control cost in the linear case. We show that the attack can increase the cost by $8500\%$ on average with ***bounded time-series perturbations that are very miniscule to the human eye*** on real electricity demand patterns. Then, we formulate attacks on ***any differentiable target function*** with a general convex controller and approximate the optimal attack. We validate our methods on synthetic ARIMA and real world electricity demand patterns with two target functions, namely the maximum action value and absolute sum of actions.

In future work, we plan to discuss how to detect the presence of adversarial attacks, as it may not be reflected in increased control costs when attackers target other system metrics. Also, more complicated gradient descent or ascent methods with momentum, like the optimizer ADAM [21], can be used in Eq. 12 to avoid local extrema. However, we need more research to analytically characterize the effect of such methods. Finally, we plan to generate certified defenses for adversarial attacks on timeseries. That is, when we know an adversarial source is perturbing the timeseries, how can a controller adjust its optimization parameters to retain its original control cost and actions?

## References

[1] S. Chinchali, P. Hu, T. Chu, M. Sharma, M. Bansal, R. Misra, M. Pavone, and S. Katti, "Cellular network traffic scheduling with deep reinforcement learning," in *Thirty-second AAAI conference on artificial intelligence*, 2018.

[2] J. Cheng, M. Pavone, S. Katti, S. Chinchali, and A. Tang, "Data sharing and compression for cooperative networked control," *Advances in Neural Information Processing Systems*, vol. 34, pp. 5947–5958, 2021.

[3] P. Donti, B. Amos, and J. Z. Kolter, "Task-based end-to-end model learning in stochastic optimization," in *Advances in Neural Information Processing Systems*, vol. 30, Curran Associates, Inc., 2017.

[4] R. Dang-Nhu, G. Singh, P. Bielik, and M. Vechev, "Adversarial attacks on probabilistic autoregressive forecasting models," in *International Conference on Machine Learning*, pp. 2356–2365, PMLR, 2020.

[5] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, "Adversarial examples are not bugs, they are features," *Advances in neural information processing systems*, vol. 32, 2019.

[6] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. El Ghaoui, and M. Jordan, "Theoretically principled trade-off between robustness and accuracy," in *International conference on machine learning*, pp. 7472–7482, PMLR, 2019.

[7] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta, "Robust adversarial reinforcement learning," in *Proceedings of the 34th International Conference on Machine Learning*, vol. 70 of *Proceedings of Machine Learning Research*, pp. 2817–2826, PMLR, 06–11 Aug 2017.

[8] X. Ma, K. Driggs-Campbell, and M. J. Kochenderfer, "Improved robustness and safety for autonomous vehicle control with adversarial reinforcement learning," in *2018 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1665–1671, 2018.

[9] A. Gleave, M. Dennis, C. Wild, N. Kant, S. Levine, and S. Russell, "Adversarial policies: Attacking deep reinforcement learning," *arXiv preprint arXiv:1905.10615*, 2019.

[10] K. Zhang, B. Hu, and T. Basar, "On the stability and convergence of robust adversarial reinforcement learning: A case study on linear quadratic systems," in *Advances in Neural Information Processing Systems*, vol. 33, pp. 22056–22068, Curran Associates, Inc., 2020.

[11] X. Zhang, X. Zhu, and L. Lessard, "Online data poisoning attacks," in *Learning for Dynamics and Control*, pp. 201–210, PMLR, 2020.

[12] A. Russo, M. Molinari, and A. Proutiere, "Data-driven control and data-poisoning attacks in buildings: the kth live-in lab case study," in *2021 29th Mediterranean Conference on Control and Automation (MED)*, pp. 53–58, 2021.

[13] U. Ghai, D. Snyder, A. Majumdar, and E. Hazan, "Generating adversarial disturbances for controller verification," in *Learning for Dynamics and Control*, pp. 1192–1204, PMLR, 2021.

[14] N. Agarwal, B. Bullins, E. Hazan, S. Kakade, and K. Singh, "Online control with adversarial disturbances," in *Proceedings of the 36th International Conference on Machine Learning* (K. Chaudhuri and R. Salakhutdinov, eds.), vol. 97 of *Proceedings of Machine Learning Research*, pp. 111–119, PMLR, 09–15 Jun 2019.

[15] A. Fu, B. Narasimhan, and S. Boyd, "CVXR: An R package for disciplined convex optimization," *Journal of Statistical Software*, vol. 94, no. 14, pp. 1–34, 2020.

[16] A. Agrawal, R. Verschueren, S. Diamond, and S. Boyd, "A rewriting system for convex optimization problems," *Journal of Control and Decision*, vol. 5, no. 1, pp. 42–60, 2018.

[17] I. I. Cplex, "V12. 1: User's manual for cplex," *International Business Machines Corporation*, vol. 46, no. 53, p. 157, 2009.

[18] M. ApS, *MOSEK Optimization Toolbox for MATLAB 9.0.105*, 2022.

[19] B. Amos and J. Z. Kolter, "Optnet: Differentiable optimization as a layer in neural networks," in *International Conference on Machine Learning*, pp. 136–145, PMLR, 2017.

[20] A. Agrawal, B. Amos, S. Barratt, S. Boyd, S. Diamond, and J. Z. Kolter, "Differentiable convex optimization layers," in *Advances in Neural Information Processing Systems*, pp. 9558–9570, 2019.

[21] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations*, 12 2014.

[22] A. C. Harvey, *ARIMA Models*, pp. 22–24. London: Palgrave Macmillan UK, 1990.

[23] "Hourly energy consumption over 10 years of hourly energy consumption data from pjm in megawatts." https://www.kaggle.com/datasets/robikscube/hourly-energy-consumption, 2018. [Online; accessed 21-June-2021].

[24] W. J. Conover, *Practical nonparametric statistics*, vol. 350. john wiley & sons, 1999.