# Adversarial Destabilization Attacks to Direct Data-Driven Control

Hampei Sasahara, *Member, IEEE,*

*Abstract*—This study investigates the vulnerability of direct data-driven control methods, specifically for the linear quadratic regulator problem, to adversarial perturbations in collected data used for controller synthesis. We consider stealthy attacks that subtly manipulate offline-collected data to destabilize the resulting closed-loop system while evading detection. To generate such perturbations, we propose the Directed Gradient Sign Method (DGSM) and its iterative variant (I-DGSM), adaptations of the fast gradient sign method originally developed for neural networks, which align perturbations with the gradient of the spectral radius of the closed-loop matrix to reduce stability. A key contribution is an efficient gradient computation technique based on implicit differentiation through the Karush-Kuhn-Tucker conditions of the underlying semidefinite program, enabling scalable and exact gradient evaluation without repeated optimization computations. To defend against these attacks, we propose two defense strategies: a regularization-based approach that enhances robustness by suppressing controller sensitivity to data perturbations and a robust data-driven control approach that guarantees closed-loop stability within bounded perturbation sets. Extensive numerical experiments on benchmark systems show that adversarial perturbations with magnitudes up to ten times smaller than random noise can destabilize controllers trained on corrupted data and that the proposed defense strategies effectively mitigate attack success rates while maintaining control performance. Additionally, we evaluate attack transferability under partial knowledge scenarios, highlighting the practical importance of protecting training data confidentiality. This work advances understanding of security risks in data-driven control and provides both attack methodologies and principled defenses critical for deploying reliable controllers in safety-critical cyber-physical systems.

*Index Terms*—Adversarial attacks, control system security, direct data-driven control.

## I. INTRODUCTION

CYBER-PHYSICAL systems (CPS), such as autonomous vehicles, consist of multiple interconnected layers, among which the perception and decision-making layer and the control layer are fundamental [1]. The perception and decision-making layer handles high-level interpretation, planning, and strategic decisions based on processed data, often leveraging artificial intelligence and machine learning techniques. In contrast, the control layer is responsible for low-level closed-loop feedback control that directly governs physical system dynamics to ensure stability and performance. It is well established that the perception and decision-making layer is vulnerable to various adversarial attacks that manipulate input data, leading to incorrect interpretations and potentially unsafe outcomes [2]–[9]. By contrast, the control layer has

traditionally been viewed as more secure due to its reliance on well-defined physical models.

However, recent advances in data-driven control paradigms have significantly shifted this perspective [10]–[12]. Data-driven control methods design controllers directly from observed input-output data without explicit system identification, enabling flexible and adaptive control even when precise models are unavailable [13]–[15]. This growing reliance on data at the control layer introduces new attack surfaces: if an adversary can manipulate the data used for controller synthesis, they may induce instability or degrade performance in ways previously unconsidered. Consequently, the vulnerability of the control layer to data manipulation has become a critical concern. Despite its importance, this area remains underexplored compared to the perception layer's adversarial risks, with only a few recent studies addressing this issue [16]–[18]. Addressing vulnerabilities in both layers, especially focusing on the emerging threats posed by data-driven control, is essential for ensuring the comprehensive security and resilience of modern CPS.

In this study, we investigate the vulnerability of direct data-driven control [11] to adversarial perturbations. Specifically, we consider an attacker who seeks to destabilize a closed-loop system by making small, intentional modifications to the input-output data used for controller synthesis. We focus on the linear quadratic regulator (LQR) problem, a widely used standard benchmark control algorithm, formulated here using a direct data-driven approach expressed as a semidefinite program (SDP) [19]. To assess worst-case risks, we assume a powerful white-box adversary with full knowledge of the system dynamics, the controller design algorithm, and the clean data. As a specific adversarial strategy, we propose the directed gradient sign method (DGSM), an adaptation of the fast gradient sign method (FGSM) originally developed for attacking neural networks [3]. DGSM computes perturbations aligned with the gradient of the spectral radius (or eigenvalues) of the closed-loop system matrix, aiming to reduce system stability. We further extend this approach with the iterative DGSM (I-DGSM), formulated as a projected gradient descent method, which refines the perturbations through multiple steps to increase their effectiveness.

A key challenge in applying I-DGSM lies in computing the gradient of the spectral radius with respect to the data. Standard numerical differentiation, such as the central difference method, would require solving an SDP multiple times proportional to the data dimension, leading to significant computational costs. To overcome this, we develop a novel approach based on implicit differentiation through the Karush-Kuhn-Tucker (KKT) conditions of the SDP, which requires solving the semidefinite program only once. This leads to exact gradient computation at a significantly reduced cost, enabling

Hampei Sasahara is with Department of Systems and Control Engineering, Institute of Science Tokyo, Tokyo 152-8552, Japan (e-mail: sasahara@sc.eng.isct.ac.jp).

efficient generation of adversarial perturbations.

Furthermore, in response to these vulnerabilities, we propose two defense methods: a regularization-based approach and a robust data-driven control approach. Although the regularization has been originally introduced to mitigate random disturbances [20], [21], we show both theoretically and numerically that it is also effective against adversarial attacks. The robust data-driven control approach builds on the recent work [22], [23] and guarantees closed-loop stability under bounded perturbations defined by a matrix ellipsoid. This robustness framework provides theoretical assurances critical for deploying data-driven controllers in safety-critical settings.

To validate the severity of adversarial risks, we conduct extensive simulations on benchmark linear time-invariant systems, including five standard mechanical and aerospace models from the tutorials on control engineering [24], [25], as well as a triple-tank system representing a typical process control application [26], [27]. A representative example with the triple tank system shows that perturbations with magnitudes as small as 0.1% to 0.5% of the data amplitude can induce instability while remaining visually indistinguishable from clean data. Across all tested systems, our results consistently demonstrate that adversarial perturbations with magnitudes approximately one-tenth or less than those of uniformly distributed random noise are sufficient to destabilize the closed-loop system. This highlights the precision and impact of the proposed adversarial perturbation algorithms in identifying minimal yet highly disruptive data manipulations. Furthermore, the proposed implicit differentiation method improves computational efficiency dramatically: compared to numerical differentiation, it achieves speedups of around 70 times, with further gains, up to nearly 3 times, when analytical derivatives are used. This efficiency enables practical evaluation of iterative attack algorithms. In addition, the numerical experiments demonstrate that the proposed defense strategies significantly reduce the attack success rate. Specifically, the regularization-based approach decreases the attack success rate from near 100 % to as low as 5% or below in several benchmark systems while maintaining control performance close to the nominal optimal. Meanwhile, the robust data-driven control method guarantees closed-loop stability under bounded perturbations and achieves nominal performance for small perturbation levels, although its control performance degrades as the allowable perturbation size increases. Finally, beyond worst-case white-box attacks, we explore effectiveness in more realistic adversarial scenarios involving limited attacker knowledge, referred to as data transferability. Specifically, we consider gray-box attacks where the adversary lacks knowledge of the exact data. Numerical results indicate that data transferability is moderate under the threat model and suggest that maintaining the confidentiality of the training data can serve as a practical defense against adversarial perturbations in control systems. This investigation provides insights into the resilience of data-driven control under less idealized but practically relevant threat models.

This work makes the following key contributions. First, we reveal a fundamental vulnerability of direct data-driven LQR control methods to adversarial perturbations in the training data. Second, we develop efficient perturbation algorithms, DGSM and its iterative variant I-DGSM, tailored to destabilize the resulting closed-loop system by exploiting the structure of the underlying SDP. Third, we introduce an efficient attack synthesis approach using implicit differentiation through the KKT conditions of the control design problem. Fourth, we propose two proactive defense strategies: a regularization-based method and a robust data-driven control approach. Finally, we validate our methods through extensive numerical experiments demonstrating both the severity of the attacks and the effectiveness of the proposed defenses, including transferability analysis under partial knowledge scenarios.

Preliminary versions of this work have been presented in [28] and [29]. In [28], the DGSM algorithm, a fundamental method for generating adversarial perturbations, was introduced. An extension of this method, referred to as I-DGSM, was proposed in [29]. Both prior works rely on numerical differentiation to compute the gradients required by these algorithms. In contrast, the present work proposes a more efficient gradient computation method, as detailed in Sec. IV. Furthermore, in addition to the regularization-based defense strategy considered in the previous papers, we introduce a new defense method based on robust data-driven control in Sec. V. Lastly, our numerical evaluation in Sec. VI includes a broader and more diverse set of benchmark systems, providing a more comprehensive validation of the proposed methods.

*Organization and Notation*

The paper is organized as follows. Sec. II provides the necessary background on direct data-driven LQR control and a common adversarial attack method against neural networks. Sec. III introduces the threat model, describing the attacker's objective and capabilities, and presents our proposed perturbation generation algorithms, including DGSM and its iterative variant. Sec. IV presents an efficient gradient computation method via implicit differentiation using the KKT conditions. Sec. V introduces defense strategies based on regularization and robust data-driven controller design. Sec. VI presents numerical experiments on benchmark systems to visualize adversarial attacks, evaluate their impact, assess the computational efficiency of the proposed gradient computation method, demonstrate the effectiveness of the proposed defense, and examine the transferability of adversarial attacks across datasets. Finally, Sec. VII concludes the paper and discusses future research directions. Appendix summarizes the basic rules of matrix calculus and exhibits proofs of the propositions.

We denote the set of real numbers by $\mathbb{R}$, the $n$-dimensional Euclidean space by $\mathbb{R}^n$, the $n$-dimensional identity matrix by $I_n$, the $n \times m$-dimensional zero matrix by $0_{n,m}$, the transpose of a matrix $M$ by $M^\mathsf{T}$, the trace and the spectral radius of a square matrix $M$ by $\mathrm{tr}(M)$ and $\rho(M)$, respectively, the 2-induced norm of a matrix $M$ by $\|M\|_2$, the Frobenius norm of a matrix $M$ by $\|M\|_\mathrm{F}$, the element-wise max norm of a matrix $M$ by $\|M\|_\mathrm{max}$, the minimum singular value of a matrix $M$ by $\sigma_\mathrm{min}(M)$, the Moore-Penrose pseudoinverse of a matrix $M$ by $M^\dagger$, the positive and negative (semi)definiteness of a Hermetian matrix $M$ by $M \succ (\succeq) 0$ and $M \prec (\preceq) 0$, respectively, the block diagonal matrix whose diagonal blocks

are composed of $M_1$ and $M_2$ by $\mathrm{diag}(M_1, M_2)$, the Kronecker product by $\otimes$, the $(p, q)$ commutation matrix by $C_{p,q}$, the component-wise sign function by $\mathrm{sign}(\cdot)$, the vectorization operator by $\mathrm{vec}$, and the Frobenius inner product of two real matrices $(M, N)$ by $\langle M, N \rangle_\mathrm{F} := \mathrm{tr}(M^\mathsf{T} N)$. The subscript for the dimension is often omitted when it is clear from the context.

## II. PRELIMINARIES

### A. Direct Data-driven LQR Control

Consider a discrete-time linear time-invariant (LTI) system

$$x_{t+1} = Ax_t + Bu_t, \quad t = 0, 1, \ldots$$

where $x_t \in \mathbb{R}^n$ is the state and $u_t \in \mathbb{R}^m$ is the control input. We assume that the pair $(A, B)$ is stabilizable. We consider the LQR problem [30, Chap. 6], which has been widely studied as a benchmark problem. Specifically, design a static state-feedback control $u_t = Kx_t$ that minimizes the cost function

$$\mathcal{J}(K) = \sum_{i=1}^n \sum_{t=0}^\infty \left( x_t^\mathsf{T} Q x_t + u_t^\mathsf{T} R u_t \right) |_{x_0=e_i}$$

with $Q \succeq 0$ and $R \succ 0$ where $e_i$ is the $i$th canonical basis vector. It is known that the cost function can be rewritten as

$$\mathcal{J}(K) = \mathrm{tr}(QP) + \mathrm{tr}(K^\mathsf{T} RKP)$$

where $P \succeq I$ is the controllability Gramian of the closed-loop system when $A + BK$ is Schur, i.e., $\rho(A + BK) < 1$. The objective of direct data-driven control is to design the optimal feedback gain using data of input and state signals without requiring explicit system identification.

The overall implementation of the direct data-driven LQR approach is structured into two distinct phases: an *offline design phase* and an *online operation phase*, which are analogous to the training and test phases commonly seen in machine learning. In the offline design phase, historical trajectory data of the system are collected under open-loop or exploratory input signals. Using this dataset, the controller is synthesized by solving an optimization problem, which learns the optimal feedback gain. This design process is entirely performed offline and does not require explicit knowledge of the system matrices $(A, B)$. Once the controller gain $K$ is computed, the online phase involves real-time operation of the system using the static linear state-feedback law $u_t = Kx_t$.

Here, we review the offline design phase in detail. It is assumed that the system matrices $(A, B)$ are unknown, but instead, the tuple of $T$-long offline batch data $(Z, X, U)$ that obey the dynamics

$$Z = AX + BU,$$

where $Z \in \mathbb{R}^{n \times T}$, $X \in \mathbb{R}^{n \times T}$, and $U \in \mathbb{R}^{m \times T}$, are available. We denote the collective data by $D := [Z^\mathsf{T} \ X^\mathsf{T} \ U^\mathsf{T}]^\mathsf{T}$. We assume that $\Gamma := [U^\mathsf{T} \ X^\mathsf{T}]^\mathsf{T}$ is full rank, which is generally necessary for data-driven LQR design [31]. It is satisfied if the input signal is persistently exciting as shown by the Willems' fundamental lemma [32].

The key idea of the approach laid out in [19] is to parameterize the controller using the available data by introducing a new variable $G \in \mathbb{R}^{T \times n}$ with the relationship

$$[K^\mathsf{T} \ I]^\mathsf{T} = \Gamma G. \tag{1}$$

Then the closed-loop matrix can be parameterized directly by data matrices as

$$A + BK = [B \ A]\Gamma G = ZG. \tag{2}$$

The LQR controller design can be formulated as

$$\begin{aligned} \min_{P,K,G} \quad & \mathcal{J}(K) \\ \text{s.t.} \quad & ZGPG^\mathsf{T} Z^\mathsf{T} - P + I \preceq 0, \\ & P \succeq I \text{ and (1)}. \end{aligned} \tag{3}$$

However, it has been pointed out that the formulation (3) is not robust to disturbance [20]. Specifically, the study considers a control system subject to disturbances, which leads to the data relationship $Z = AX + BU + W$ where $W \in \mathbb{R}^{n \times T}$ is the disturbance matrix. Then we have $A + BK = (Z - W)G$ instead of (2), and hence the optimization problem (3) fails to produce the exact solution to the LQR problem. To address this issue, a regularized formulation has been proposed for enhancing robustness against disturbance:

$$\begin{aligned} \min_{P,K,G} \quad & \mathcal{J}(K) + \gamma\|\Pi G\|_\mathrm{F}^2 \\ \text{s.t.} \quad & ZGPG^\mathsf{T} Z^\mathsf{T} - P + I \preceq 0, \\ & P \succeq I \text{ and (1)} \end{aligned} \tag{4}$$

with a constant $\gamma \geq 0$ where $\Pi := I - \Gamma^\dagger \Gamma$. The regularizer $\gamma\|\Pi G\|_\mathrm{F}$ is referred to as certainty-equivalence regularization because it leads to the controller equivalent to the certainty-equivalence indirect data-driven LQR with ordinary least-square estimation of the system model when $\gamma$ is sufficiently large [20]. Note that the reformulated problem can also be converted into the following SDP:

$$\begin{aligned} \min_{L,S} \quad & J(L, S, D) \\ \text{s.t.} \quad & F(L, S, D) \succeq 0 \end{aligned} \tag{5}$$

with variables $L \in \mathbb{R}^{T \times n}, S \in \mathbb{R}^{m \times m}$ for given data $D$ where

$$J := \mathrm{tr}(QXL) + \mathrm{tr}(S) + \gamma\|\Pi L\|_\mathrm{F}^2, \ F := \mathrm{diag}(F_1, F_2),$$

and

$$F_1 := \begin{bmatrix} S & VUL \\ * & XL \end{bmatrix}, F_2 := \begin{bmatrix} XL - I & ZL \\ * & XL \end{bmatrix},$$

and $V \succ 0$ is the square root of $R$. Note that $F(L, S, D) \in \mathbb{R}^{(3n+m) \times (3n+m)}$. The resulting controller is given by $K = F_K(L, D)$ with $F_K(L, D) := UL(XL)^{-1}$ and $P = XL$ where $(L, S)$ is the optimal solution to (5).

### B. Fast Gradient Sign Method

The Fast Gradient Sign Method (FGSM) is a widely used technique for efficiently generating adversarial perturbations against trained neural networks [3]. Given a loss function $\ell(D, Y; \theta)$ of the targeted neural network where $D \in \mathbb{R}^{p \times q}$ is the input data, $Y \in \mathcal{Y}$ is the true label, and $\theta$ is the trained parameter, FGSM aims to craft a perturbation $\Delta \in \mathbb{R}^{p \times q}$ that causes the classifier $f : \mathbb{R}^{p \times q} \to \mathcal{Y}$ to misclassify the input,
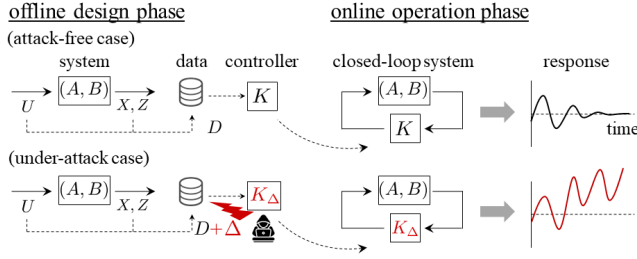
Fig. 1: Threat model addressed in this paper. The adversary is capable of manipulating the data $D$ by introducing a perturbation $\Delta$ with the knowledge of the system model, the clean data, and the controller design algorithm to be implemented. The controller $K_\Delta$ designed based on the perturbed data $D+\Delta$ may lead to instability of the closed-loop system.

i.e., $f(D + \Delta) \neq f(X)$. The perturbation is constrained in magnitude by the element-wise max norm $\|\Delta\|_{\max} \leq \epsilon$, where $\epsilon > 0$ is a small constant ensuring imperceptibility.

The core idea behind FGSM is to exploit the linear approximation of the loss function with respect to the input:

$$\ell(D + \Delta, Y; \theta) \simeq \ell(D, Y; \theta) + \langle \nabla_D \ell(D, Y; \theta), \Delta \rangle_F.$$

To maximize the loss under the norm constraint, FGSM chooses each entry of $\Delta$ to align with the sign of the gradient:

$$\Delta = \epsilon \, \mathrm{sign}(\nabla_D \ell(D, Y; \theta)).$$

This perturbation direction maximally increases the loss in a single step under the given constraint. In practice, FGSM is often used iteratively, gradually increasing $\epsilon$ until misclassification is achieved.

## III. ADVERSARIAL ATTACK FRAMEWORK

### A. Threat Model

We consider an adversary whose goal is to destabilize the closed-loop control system by introducing small but malicious perturbations to the input-state data used for controller synthesis. The adversarial manipulations are designed to render the resulting system unstable, while remaining imperceptible. Fig. 1 illustrates the overall threat model.

The attack is performed during the offline design phase. The adversary introduces a perturbation $(\Delta_Z, \Delta_X, \Delta_U)$ to the clean input and state data $(Z, X, U)$ used in the direct data-driven LQR control algorithm, thereby influencing the resulting controller learned from the perturbed data

$$(Z_\Delta, X_\Delta, U_\Delta) := (Z + \Delta_Z, X + \Delta_X, U + \Delta_U).$$

This type of attack is often referred to as data poisoning attack [33], because it corrupts the training data prior to learning, causing the learned controller to behave in a manner unintended by the designer. In this context, the "training data" corresponds to the historical trajectories used for direct controller synthesis, and the adversary exploits this dependency to inject targeted instabilities.

The adversary can modify the input and state sequences used for control design, but the magnitude of perturbation is bounded. We constrain the norm of the perturbation matrix $\Delta := [\Delta_Z^\mathsf{T} \ \Delta_X^\mathsf{T} \ \Delta_U^\mathsf{T}]^\mathsf{T}$, i.e., $\|\Delta\| \leq \epsilon$ with a small constant $\epsilon > 0$ with a chosen matrix norm $\|\cdot\|$. Specifically, we adopt the element-wise max norm $\|\cdot\|_{\max}$, because it directly limits the largest allowable change to any individual data point. Alternatively, one may consider a relative perturbation constraint, each of $(\Delta_Z, \Delta_X, \Delta_U)$ is bounded in proportion to the corresponding data matrix, i.e.,

$$\|\Delta_Z\| \leq \epsilon \|Z\|, \quad \|\Delta_X\| \leq \epsilon \|X\|, \quad \|\Delta_U\| \leq \epsilon \|U\|. \quad (6)$$

Such a constraint captures scenarios in which the adversary is further restricted to preserving the shape or scale of the original signal, thereby enhancing the stealthiness of the attack under stricter plausibility requirements.

We primarily consider a white-box adversary with full knowledge of the true system dynamics $(A, B)$, the controller design algorithm, and the clean input-state data $D$. While this represents a worst-case scenario that establishes an upper bound on potential attack impact, such a powerful adversary may be unrealistic in practice. To address this limitation, we later extend our analysis to a gray-box setting, where the adversary lacks knowledge of the precise dataset used for controller design. In this case, a plausible attack strategy is to use hypothetical data $D_{\mathrm{hyp}}$ that is consistent with the system dynamics but follows a different trajectory from the training data. We refer to the effectiveness of such an attack as transferability across data, numerically evaluated in Sec. VI.

*Remark:* Such adversarial perturbations may arise in practice through compromised sensing or communication channels during the data collection phase. For instance, if the input-state trajectories used for controller design are transmitted over a network, a man-in-the-middle attacker could subtly tamper with the data in transit [34]. Similarly, malware residing on a data logger or edge device could slightly modify recorded sensor or actuator signals before they are used in control design [35]. These scenarios are particularly concerning in industrial or cyber-physical systems, where data integrity cannot always be guaranteed and controller updates may be deployed based solely on observed data without re-verifying physical models.

### B. Perturbation Generation Algorithm

*1) Directed Gradient Sign Method (DGSM):* Having established the threat model, we now formulate the perturbation generation task as an optimization problem. The adversary's objective is to craft a perturbation to the input-state data such that the controller synthesized from the perturbed data leads to an unstable closed-loop system. This task is formalized as maximizing the spectral radius of the closed-loop matrix, subject to a norm constraint on the perturbation to ensure plausibility. Formally, let $\mathcal{K} : D \mapsto K$ denote the mapping from the given data $D$ to the controller designed by the direct data-driven control described in Sec. II. Then the optimization problem can be formulated as

$$\begin{aligned} \max_\Delta \quad & \rho(A + B\mathcal{K}(D + \Delta)) \\ \text{s.t.} \quad & \|\Delta\|_{\max} \leq \epsilon. \end{aligned}$$
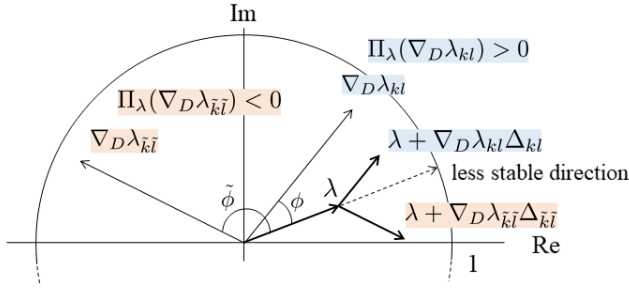
Fig. 2: Function $\Pi_\lambda$ in (7) plays a pivotal role. Due to the alignment of $\nabla_D \lambda_{kl}$ with the direction of $\lambda$, the angle $\phi$ formed between $\lambda$ and $\nabla_D \lambda_{kl}$ remains less than $\pi/2$, thereby causing $\Pi_\lambda(\nabla_D \lambda_{kl}) > 0$. Conversely, the angle $\tilde{\phi}$ formed between $\lambda$ and $\nabla_D \lambda_{\tilde{k}\tilde{l}}$ is greater than $\pi/2$, leading to $\Pi_\lambda(\nabla_D \lambda_{\tilde{k}\tilde{l}}) < 0$. As a result, in both cases, the eigenvalue subjected to the perturbation is shifted closer to the unit circle.

Our proposed algorithm, Directed Gradient Sign Method (DGSM), is a natural extension of FGSM tailored to the control-theoretic setting of our problem. DGSM generates adversarial perturbations that aim to destabilize the closed-loop system by aligning the perturbation with the gradient of the spectral radius of the closed-loop matrix. Specifically, it chooses

$$\Delta = \epsilon \operatorname{sign}(\nabla_D \rho(A + B\mathcal{K}(D))),$$

assuming that $\rho$ is locally differentiable. Geometrically, this perturbation shifts the dominant eigenvalue $\lambda(D + \Delta)$, which corresponds to the spectral radius of the closed-loop matrix, in the direction that reduces stability. Note that the linear approximation of the perturbed dominant eigenvalue is given by

$$\lambda(D + \Delta) \simeq \lambda + \langle \nabla_D \lambda, \Delta \rangle_{\mathrm{F}}.$$

The perturbation matrix computed by DGSM can be equivalently expressed as

$$\Delta_{kl} = \epsilon \operatorname{sign}(\Pi_\lambda(\nabla_D \lambda_{kl}))$$

where the subscript $kl$ denotes the $(k, l)$th component and $\Pi_\lambda : \mathbb{C} \to \mathbb{R}$ is an operator defined by

$$\Pi_\lambda(\nabla_D \lambda_{kl}) := \operatorname{Re}(\lambda)\operatorname{Re}(\nabla_D \lambda_{kl}) + \operatorname{Im}(\lambda)\operatorname{Im}(\nabla_D \lambda_{kl}). \tag{7}$$

Fig. 2 illustrates a geometric interpretation of the operator $\Pi_\lambda$. Suppose that $\nabla_D \lambda_{kl}$ points in the direction of $\lambda$. More precisely, the angle between $\lambda$ and $\nabla_D \lambda_{kl}$, denoted by $\phi$, is less than $\pi/2$, which leads to $\Pi_\lambda(\nabla_D \lambda_{kl}) > 0$. We now suppose that the angle between $\lambda$ and another element $\nabla_D \lambda_{kl}$, denoted by $\tilde{\phi}$, is greater than $\pi/2$. Then we have $\Pi_\lambda(\nabla_D \lambda_{\tilde{k}\tilde{l}}) < 0$. In both cases, owing to the function $\Pi_\lambda$, the perturbed eigenvalue moves closer to the unit circle. This alignment with the destabilizing direction motivates the name "Directed" Gradient Sign Method.

*2) Iterative Directed Gradient Sign Method (I-DGSM):* To enhance the performance of DGSM, we also propose an iterative variant based on DGSM, referred to as Iterative Directed Gradient Sign Method (I-DGSM). Algorithm 1 provides its

algorithm where $\alpha_{\mathrm{step}} > 0$ is the step size and $\operatorname{Proj}_\epsilon(\hat{\Delta})$ denotes the element-wise truncation function, i.e., its $(k, l)$th component is given by

$$\operatorname{Proj}_\epsilon(\hat{\Delta})_{kl} = \begin{cases} \epsilon \operatorname{sign}(\hat{\Delta}_{kl}) & \text{if } |\hat{\Delta}_{kl}| \geq \epsilon, \\ \hat{\Delta}_{kl} & \text{otherwise.} \end{cases}$$

In the relative constraint case with (6), the projection map is individually applied to $(\Delta_Z, \Delta_X, \Delta_U)$ with the corresponding range. This approach adopts a projected gradient ascent framework [36, Chap. 3] to more effectively solve the underlying perturbation generation problem. Recall that the original task of crafting adversarial perturbations can be formulated as a bi-level optimization problem, where the inner problem involves synthesizing a controller from perturbed data, and the outer problem seeks to maximize the spectral radius of the resulting closed-loop system. I-DGSM addresses this by iteratively updating the perturbation in the direction of the gradient of the spectral radius and enforcing the perturbation norm constraint via projection after each step. While global optimality cannot be guaranteed due to the non-convex nature of the bi-level problem, the algorithm is designed to converge to a local maximum with a sufficiently small step size $\alpha_{\mathrm{step}}$, providing a practical algorithm to identify severe, stability-compromising perturbations. Note that, although this algorithm uses vanilla gradient ascent with a fixed step size for simplification, we can use another gradient-based optimization method, such as the adaptive gradient algorithm (Adagrad) [37].

---

**Algorithm 1** Iterative Directed Gradient Sign Method (I-DGSM)

---

**Input:** $\epsilon, \alpha_{\mathrm{step}}, A, B, D, \mathcal{K}$
**Output:** $\Delta$
1: $\Delta_0 = 0$
2: $k \leftarrow 0$
3: **while** Termination condition not met **do**
4: $\quad \hat{\Delta}_{k+1} \leftarrow \Delta_k + \alpha_{\mathrm{step}} \nabla_D \rho(A + B\mathcal{K}(D + \Delta_k))$
5: $\quad \Delta_{k+1} \leftarrow \operatorname{Proj}_\epsilon(\hat{\Delta}_{k+1})$
6: $\quad k \leftarrow k + 1$
7: **end while**
8: **return** $\Delta_k$

---

## IV. EFFICIENT GRADIENT COMPUTATION

To implement I-DGSM, it is essential to compute the gradient of the spectral radius of the closed-loop matrix with respect to the perturbation matrix. The chain rule (see Appendix A) leads to the relationship

$$\frac{d\rho}{dD} = \frac{d\rho}{d\mathcal{K}}\left(\frac{\partial F_K}{\partial L}\frac{dL}{dD} + \frac{\partial F_K}{\partial D}\right) \tag{8}$$

Here, $d\rho/d\mathcal{K}$, $\partial F_K/\partial L$, and $\partial F_K/\partial D$ can be efficiently computed numerically. In contrast, the computation of $dL/dD$ is non-trivial because the controller is not given in a closed form, and instead, it is the solution to an SDP whose constraints and objective depend on the perturbed data. As a result, the closed-loop matrix becomes an implicit function of the perturbation, complicating the gradient computation. A

straightforward approach based on numerical differentiation (e.g., central difference method [38, Chap. 4]) requires solving the SDP $\mathcal{O}((2n + m)T)$ times. This requirement is computationally expensive and impractical, especially for I-DGSM where the gradient must be computed at every iteration. This motivates the need for a more efficient and scalable method to compute the required gradients.

### A. Implicit Differentiation Approach

To address the computational overhead of numerical differentiation, we adopt an *implicit differentiation* approach [39]–[41], by leveraging the optimality conditions of the optimization problem used in controller synthesis. Instead of explicitly solving perturbed SDPs repeatedly to estimate the gradient, we derive an analytical expression for the gradient by differentiating an implicit function defined by $F_{\text{imp}}(L(D), D) = 0$, where we express the argument $D$ for the optimal solution $L$ to emphasize the dependency. Applying the chain rule yields

$$\frac{\partial F_{\text{imp}}}{\partial L} \frac{dL}{dD} + \frac{\partial F_{\text{imp}}}{\partial D} = 0, \quad (9)$$

from which $dL/dD$ can be obtained by solving the linear equation, provided the partial derivatives $\partial F_{\text{imp}}/\partial L$ and $\partial F_{\text{imp}}/\partial D$ are tractable.

To construct the implicit function, we exploit the KKT condition for SDP [42, Chap. 5]. Specifically, the KKT condition for the SDP (5) is given as

$$\begin{cases} \partial \mathcal{L}(L, S, \Lambda, D)/\partial L = 0, \\ \partial \mathcal{L}(L, S, \Lambda, D)/\partial S = 0, \\ \text{tr}(F(L, S, D)\Lambda^{\mathsf{T}}) = 0, \\ \Lambda - \Lambda^{\mathsf{T}} = 0, \end{cases} \quad (10)$$

with the primal and dual feasibility conditions $F(L, S, D) \succeq 0$ and $\Lambda \succeq 0$, where $\Lambda \in \mathbb{R}^{(3n+m) \times (3n+m)}$ denotes the Lagrange multipliers and the Lagrangian $\mathcal{L}$ is given by $\mathcal{L}(L, S, \Lambda, D) := J(L, S, D) - \text{tr}(F(L, S, D)\Lambda^{\mathsf{T}})$. In our approach, we focus on the equality constraints in (10), excluding the semidefiniteness constraints. This relaxation is justified by prior work [39], [43], which show that the primal and dual feasibility conditions can be neglected around the optimal solution. We hereinafter assume the existence of a continuously differentiable implicit function $(L(D), S(D), \Lambda(D))$ locally around the solution satisfying the KKT conditions, which can be verified by applying the implicit function theorem [44, Chap. 7]. For notational simplicity, we define the four components of the KKT conditions as $G_i = 0$ for $i = 1, \ldots, 4$, where $G_1 \in \mathbb{R}^{1 \times nT}$, $G_2 \in \mathbb{R}^{1 \times m^2}$, $G_3 \in \mathbb{R}$, $G_4 \in \mathbb{R}^{(3n+m) \times (3n+m)}$.

Differentiating the system (10) with respect to $D$ yields the linear equation

$$\underbrace{\begin{bmatrix} \dfrac{\partial G_1}{\partial L} & \dfrac{\partial G_1}{\partial S} & \dfrac{\partial G_1}{\partial \Lambda} \\ \vdots & \vdots & \vdots \\ \dfrac{\partial G_4}{\partial L} & \dfrac{\partial G_4}{\partial S} & \dfrac{\partial G_4}{\partial \Lambda} \end{bmatrix}}_{:=H} \begin{bmatrix} \dfrac{dL}{dD} \\ \dfrac{dS}{dD} \\ \dfrac{d\Lambda}{dD} \end{bmatrix} + \begin{bmatrix} \dfrac{\partial G_1}{\partial D} \\ \vdots \\ \dfrac{\partial G_4}{\partial D} \end{bmatrix} = 0, \quad (11)$$

which corresponds to (9). Note that the coefficient matrices are comparatively easy to compute numerically because we have closed-form expressions for $\mathcal{L}, G_3, G_4$.

### B. Rank Complement

The linear system (11), derived from differentiating the KKT conditions, can be written compactly using a coefficient matrix $H \in \mathbb{R}^{(nT+m^2+1+(3n+m)^2) \times (nT+m^2+(3n+m)^2)}$. However, this equation is generally underdetermined due to rank deficiency of $H$:

**Proposition 1** The matrix $H$ has a non-trivial kernel for any $L, S, \Lambda, D$.

This rank deficiency arises from the structure of the KKT conditions, in particular the scalar complementarity condition $\text{tr}(F\Lambda^{\mathsf{T}}) = 0$. While this constraint is sufficient to certify optimality, it introduces only a single scalar equation, which limits the rank of $H$ and prevents unique identification of the gradient $dL/dD$. To resolve this issue, we seek to inject additional valid constraints that are consistent with optimality but provide more structural information. To this end, we leverage the following result:

**Proposition 2** For any positive semidefinite matrices $S_1$ and $S_2$, $\text{tr}(S_1 S_2) = 0$ if and only if $S_1 S_2 = 0$.

Proposition 2 implies that, under the positive semidefinite assumption, the scalar condition $G_3 = \text{tr}(F\Lambda^{\mathsf{T}}) = 0$ can be equivalently replaced by the matrix equality $\bar{G}_3 := F\Lambda^{\mathsf{T}} = 0$. This matrix equality introduces multiple scalar equations and thus contributes more rank to the system. By incorporating this into the KKT structure, we obtain the following augmented system:

$$\begin{bmatrix} \dfrac{\partial G_1}{\partial L} & \dfrac{\partial G_1}{\partial S} & \dfrac{\partial G_1}{\partial \Lambda} \\ \dfrac{\partial G_2}{\partial L} & \dfrac{\partial G_2}{\partial S} & \dfrac{\partial G_2}{\partial \Lambda} \\ \dfrac{\partial \bar{G}_3}{\partial L} & \dfrac{\partial \bar{G}_3}{\partial S} & \dfrac{\partial \bar{G}_3}{\partial \Lambda} \\ \dfrac{\partial G_4}{\partial L} & \dfrac{\partial G_4}{\partial S} & \dfrac{\partial G_4}{\partial \Lambda} \end{bmatrix} \begin{bmatrix} \dfrac{dL}{dD} \\ \dfrac{dS}{dD} \\ \dfrac{d\Lambda}{dD} \end{bmatrix} + \begin{bmatrix} \dfrac{\partial G_1}{\partial D} \\ \dfrac{\partial G_2}{\partial D} \\ \dfrac{\partial \bar{G}_3}{\partial D} \\ \dfrac{\partial G_4}{\partial D} \end{bmatrix} = 0. \quad (12)$$

The complete gradient computation procedure is summarized in Algorithm 2. Numerical examples in Sec. VI empirically indicate the uniqueness of the solution to the augmented system (12).

---

**Algorithm 2** Gradient Computation using Implicit Differentiation

---

**Input:** $J, F, D$
**Output:** $dL/dD$
 1: Solve (5) given $D$ and obtain optimal solutions $(L, S, \Lambda)$
 2: Compute the coefficient matrices in (12) around $(L, S, \Lambda)$
 3: Solve (12) and obtain $(dL/dD, dS/dD, d\Lambda/dD)$
 4: **return** $dL/dD$

---

## C. Further Efficient Computation

Computing the coefficient matrices in (12) via numerical differentiation requires evaluating $\mathcal{L}$ many times. This leads to a significant computational burden and can become a bottleneck for calculation precision, especially for large-scale problems. To address this, we derive closed-form expressions for the required derivatives, which substantially reduce computational cost. Throughout the following analysis, we assume differentiability of the functions. We begin by presenting the analytical derivatives used in (8).

**Proposition 3** We have

$$d\rho/dK = \mathrm{Re}\left[\lambda \,\overline{w^{\mathsf{T}}(v^{\mathsf{T}} \otimes B)/(w^{\mathsf{T}}v)}\right]/|\lambda|,$$
$$\partial F_K/\partial L = (XL)^{-\mathsf{T}} \otimes (U - F_K X),$$
$$\partial F_K/\partial D = ((XL)^{-\mathsf{T}} L^{\mathsf{T}}) \otimes (E_U - F_K E_X)$$

where the overline denotes the complex conjugate, $\lambda, v, w$ are the dominant eigenvalue of $A + BK$ and the corresponding right and left eigenvectors, respectively, and $E_U := [0_{m,n}\ 0_{m,n}\ I_m]$ and $E_X := [0_{n,n}\ I_n\ 0_{n,m}]$.

Subsequently, we derive the derivative of the Lagrangian.

**Proposition 4** We have

$$\frac{\partial \mathcal{L}}{\partial L} = (\mathrm{vec}(X^{\mathsf{T}} Q^{\mathsf{T}}))^{\mathsf{T}} + 2\gamma(\mathrm{vec}(\Pi L))^{\mathsf{T}} - (\mathrm{vec}(\Lambda))^{\mathsf{T}}\frac{\partial F}{\partial L},$$
$$\frac{\partial \mathcal{L}}{\partial S} = (\mathrm{vec}(I_m))^{\mathsf{T}} - (\mathrm{vec}(\Lambda))^{\mathsf{T}}\frac{\partial F}{\partial S}, \tag{13}$$

with

$$\partial F/\partial L = E_1 \partial F_1/\partial L + E_2 \partial F_2/\partial L,$$
$$\partial F/\partial S = E_1 \partial F_1/\partial S + E_2 \partial F_2/\partial S, \tag{14}$$

where $E_1 := \bar{E}_1 \otimes \bar{E}_1$, $E_2 := \bar{E}_2 \otimes \bar{E}_2$,

$$\bar{E}_1 := \begin{bmatrix} I_{n+m} \\ 0_{2n,n+m} \end{bmatrix}, \quad \bar{E}_2 := \begin{bmatrix} 0_{n+m,2n} \\ I_{2n} \end{bmatrix},$$

and the derivatives of $F_1$ and $F_2$ are given in (15).

Finally, we obtain the coefficient matrices in (12).

**Proposition 5** We have

$$\frac{\partial G_1}{\partial L} = 2\gamma(I_n \otimes \Pi), \ \frac{\partial G_1}{\partial S} = 0_{nT,m^2}, \ \frac{\partial G_1}{\partial \Lambda} = -\frac{\partial F}{\partial L}^{\mathsf{T}},$$
$$\frac{\partial G_1}{\partial D} = C_{n,T}(I_T \otimes (QE_X)) + 2\gamma(L^{\mathsf{T}} \otimes I_T)\frac{d\Pi}{dD}$$
$$\qquad - (I_{nT} \otimes (\mathrm{vec}(\Lambda))^{\mathsf{T}})\frac{\partial^2 F}{\partial D \partial L},$$
$$\frac{\partial G_2}{\partial L} = 0_{m^2,nT}, \ \frac{\partial G_2}{\partial S} = 0_{m^2,m^2}, \ \frac{\partial G_2}{\partial \Lambda} = -\frac{\partial F}{\partial S}^{\mathsf{T}},$$
$$\frac{\partial G_2}{\partial D} = -(I_{m^2} \otimes (\mathrm{vec}(\Lambda))^{\mathsf{T}})\frac{\partial^2 F}{\partial D \partial S},$$
$$\frac{\partial G_3}{\partial L} = (\Lambda \otimes I_{\hat{n}})\frac{\partial F}{\partial L}, \ \frac{\partial G_3}{\partial S} = (\Lambda \otimes I_{\hat{n}})\frac{\partial F}{\partial S},$$
$$\frac{\partial G_3}{\partial \Lambda} = (I_{\hat{n}} \otimes F)C_{\hat{n},\hat{n}}, \ \frac{\partial G_3}{\partial D} = (\Lambda \otimes I_{\hat{n}})\frac{\partial F}{\partial D},$$
$$\frac{\partial G_4}{\partial L} = 0_{(\hat{n})^2,nT}, \ \frac{\partial G_4}{\partial S} = 0_{(\hat{n})^2,m^2},$$
$$\frac{\partial G_4}{\partial \Lambda} = I_{(\hat{n})^2} - C_{\hat{n},\hat{n}}, \ \frac{\partial G_4}{\partial D} = 0_{(\hat{n})^2,\bar{n}T},$$

where the derivatives are given by

$$\frac{\partial F}{\partial D} = E_1 \frac{\partial F_1}{\partial D} + E_2 \frac{\partial F_2}{\partial D},$$
$$\frac{\partial^2 F}{\partial D \partial L} = (I_{nT} \otimes E_1)\frac{\partial^2 F_1}{\partial D \partial L} + (I_{nT} \otimes E_2)\frac{\partial^2 F_2}{\partial D \partial L},$$
$$\frac{d\Pi}{dD} = \frac{\partial \Pi}{\partial X}(I_T \otimes E_X) + \frac{\partial \Pi}{\partial U}(I_T \otimes E_U),$$
$$\frac{d\Pi}{dX} = -(\Gamma \otimes I_T)\frac{d\Gamma^{\dagger}}{d\Gamma}\frac{\partial \Gamma}{\partial X} - (I_T \otimes \Gamma^{\dagger})\frac{\partial \Gamma}{\partial X},$$
$$\frac{d\Pi}{dU} = -(\Gamma \otimes I_T)\frac{d\Gamma^{\dagger}}{d\Gamma}\frac{\partial \Gamma}{\partial U} - (I_T \otimes \Gamma^{\dagger})\frac{\partial \Gamma}{\partial U},$$
$$\frac{\partial \Gamma}{\partial X} = I_T \otimes \begin{bmatrix} 0_{m,n} \\ I_n \end{bmatrix}, \quad \frac{\partial \Gamma}{\partial U} = I_T \otimes \begin{bmatrix} I_m \\ 0_{n,m} \end{bmatrix}$$

with (15), (16), $\hat{n} := 3n + m$, $\bar{n} := 2n + m$, and

$$C_1 := I_n \otimes C_{T,m+n} \otimes I_{m+n}, \quad C_2 := I_n \otimes C_{T,2n} \otimes I_{2n},$$
$$\tilde{C}_1 := I_T \otimes C_{n,m+n} \otimes I_{m+n}, \quad \tilde{C}_2 := I_T \otimes C_{n,2n} \otimes I_{2n}.$$

*Remark:* While closed-form computation of derivatives is significantly more efficient and accurate than numerical differentiation as observed in the numerical experiments, it often incurs substantial memory overhead. For instance, the second derivative of the matrix-valued function $F \in \mathbb{R}^{(3n+m) \times (3n+m)}$ with respect to $L \in \mathbb{R}^{T \times n}$ and $D \in \mathbb{R}^{(2n+m) \times T}$ can result in a dense matrix of size $(3n + m)^2 nT \times (2n + m)T$. For modest problem sizes, such as $n = m = T = 10$, the number of elements in this matrix reaches 48,000,000. Storing and manipulating such large matrices can pose serious memory and scalability challenges. Resolving this issue is included in future work.

## V. DEFENSE STRATEGIES

Countermeasures for adversarial attacks can be classified into two categories: reactive and proactive strategies [45]. While reactive strategies mainly include detection of attacks, proactive strategies include robustness enhancement of the designed architecture. In this section, we propose two proactive defense methods that enhance closed-loop stability.

### A. Stability Enhancement by Regularization

We expect that the regularization introduced in (4) can enhance stability against adversarial attacks, similar to its effect on disturbance rejection. Although its effectiveness in the context of disturbances has been theoretically analyzed in [20], this analysis does not directly extend to adversarial perturbations due to fundamental differences between them. For example, while adversarial attacks are added to the data afterward, disturbances enter the internal loop and affect all subsequent outputs. Additionally, adversarial attacks target input data, while disturbances do not alter the inputs themselves. These distinctions necessitate a separate analysis to understand how the proposed regularization contributes to stability under adversarial conditions.

Let $\mu$ be a signal-to-perturbation ratio given by

$$\mu := \min\left(\|Z\|_2/\|\Delta_Z\|_2, \sigma_{\min}(\Gamma)/\|\Delta_\Gamma\|_2\right)$$

$$\frac{\partial F_1}{\partial L} = \begin{bmatrix} 0_{m,n} \\ I_n \end{bmatrix} \otimes \begin{bmatrix} VU \\ 0_{n,T} \end{bmatrix} + \left( \begin{bmatrix} VU \\ 0_{n,T} \end{bmatrix} \otimes \begin{bmatrix} 0_{m,n} \\ I_n \end{bmatrix} \right) C_{T,n} + \begin{bmatrix} 0_{m,n} \\ I_n \end{bmatrix} \otimes \begin{bmatrix} 0_{m,T} \\ X \end{bmatrix},$$

$$\frac{\partial F_2}{\partial L} = \begin{bmatrix} I_n \\ 0_{n,n} \end{bmatrix} \otimes \begin{bmatrix} X \\ 0_{n,T} \end{bmatrix} + \begin{bmatrix} 0_{n,n} \\ I_n \end{bmatrix} \otimes \begin{bmatrix} Z \\ 0_{n,T} \end{bmatrix} + \left( \begin{bmatrix} Z \\ 0_{n,T} \end{bmatrix} \otimes \begin{bmatrix} 0_{n,n} \\ I_n \end{bmatrix} \right) C_{T,n} + \begin{bmatrix} 0_{n,n} \\ I_n \end{bmatrix} \otimes \begin{bmatrix} 0_{n,T} \\ X \end{bmatrix}, \qquad (15)$$

$$\frac{\partial F_1}{\partial S} = \begin{bmatrix} I_m \\ 0_{n,m} \end{bmatrix} \otimes \begin{bmatrix} I_m \\ 0_{n,m} \end{bmatrix}, \quad \frac{\partial F_2}{\partial S} = 0_{4n^2,m^2}.$$

---

$$\frac{\partial F_1}{\partial D} = \begin{bmatrix} 0_{m,T} \\ L^{\mathsf{T}} \end{bmatrix} \otimes \begin{bmatrix} VE_U \\ 0_{n,\bar{n}} \end{bmatrix} + \left( \begin{bmatrix} VE_U \\ 0_{n,\bar{n}} \end{bmatrix} \otimes \begin{bmatrix} 0_{m,T} \\ L^{\mathsf{T}} \end{bmatrix} \right) C_{\bar{n},T} + \begin{bmatrix} 0_{m,T} \\ L^{\mathsf{T}} \end{bmatrix} \otimes \begin{bmatrix} 0_{m,\bar{n}} \\ E_X \end{bmatrix},$$

$$\frac{\partial F_2}{\partial D} = \begin{bmatrix} L^{\mathsf{T}} \\ 0_{n,T} \end{bmatrix} \otimes \begin{bmatrix} E_X \\ 0_{n,\bar{n}} \end{bmatrix} + \begin{bmatrix} 0_{n,T} \\ L^{\mathsf{T}} \end{bmatrix} \otimes \begin{bmatrix} E_Z \\ 0_{n,\bar{n}} \end{bmatrix} + \left( \begin{bmatrix} E_Z \\ 0_{n,\bar{n}} \end{bmatrix} \otimes \begin{bmatrix} 0_{n,T} \\ L^{\mathsf{T}} \end{bmatrix} \right) C_{\bar{n},T} + \begin{bmatrix} 0_{n,T} \\ L^{\mathsf{T}} \end{bmatrix} \otimes \begin{bmatrix} 0_{n,\bar{n}} \\ E_X \end{bmatrix},$$

$$\frac{\partial^2 F_1}{\partial D \partial L} = C_1 \left( \mathrm{vec}\left( \begin{bmatrix} 0_{m,n} \\ I_n \end{bmatrix} \right) \otimes I_{(m+n)T} \right) \left( I_T \otimes \begin{bmatrix} VE_U \\ 0_{n,\bar{n}} \end{bmatrix} \right) + C_1 \left( \mathrm{vec}\left( \begin{bmatrix} 0_{m,n} \\ I_n \end{bmatrix} \right) \otimes I_{(m+n)T} \right) \left( I_T \otimes \begin{bmatrix} 0_{m,\bar{n}} \\ E_X \end{bmatrix} \right)$$

$$+ (C_{n,T}^{\mathsf{T}} \otimes I_{(m+n)^2}) \tilde{C}_1 \left( I_{(m+n)T} \otimes \mathrm{vec}\left( \begin{bmatrix} 0_{m,n} \\ I_n \end{bmatrix} \right) \right) \left( I_T \otimes \begin{bmatrix} VE_U \\ 0_{n,\bar{n}} \end{bmatrix} \right),$$

$$\frac{\partial^2 F_2}{\partial D \partial L} = C_2 \left( \mathrm{vec}\left( \begin{bmatrix} I_n \\ 0_{n,n} \end{bmatrix} \right) \otimes I_{2nT} \right) \left( I_T \otimes \begin{bmatrix} E_X \\ 0_{n,\bar{n}} \end{bmatrix} \right) + C_2 \left( \mathrm{vec}\left( \begin{bmatrix} 0_{n,n} \\ I_n \end{bmatrix} \right) \otimes I_{2nT} \right) \left( I_T \otimes \begin{bmatrix} E_Z \\ 0_{n,\bar{n}} \end{bmatrix} \right)$$

$$+ (C_{n,T}^{\mathsf{T}} \otimes I_{4n^2}) \tilde{C}_2 \left( I_{2nT} \otimes \mathrm{vec}\left( \begin{bmatrix} 0_{n,n} \\ I_n \end{bmatrix} \right) \right) \left( I_T \otimes \begin{bmatrix} E_Z \\ 0_{n,\bar{n}} \end{bmatrix} \right) + C_2 \left( \mathrm{vec}\left( \begin{bmatrix} 0_{n,n} \\ I_n \end{bmatrix} \right) \otimes I_{2nT} \right) \left( I_T \otimes \begin{bmatrix} 0_{n,\bar{n}} \\ E_X \end{bmatrix} \right),$$

$$\frac{d\Gamma^{\dagger}}{d\Gamma} = ((\Gamma\Gamma^{\mathsf{T}})^{-\mathsf{T}} \otimes I_T) C_{n+m,T} - (I_{m+n} \otimes \Gamma^{\mathsf{T}})((\Gamma\Gamma^{\mathsf{T}})^{-\mathsf{T}} \otimes (\Gamma\Gamma^{\mathsf{T}})^{-1})((\Gamma \otimes I_{m+n}) + (I_{m+n} \otimes \Gamma)) C_{m+n,T}. \qquad (16)$$

---

where $\Delta_\Gamma := [\Delta_U^{\mathsf{T}} \ \Delta_X^{\mathsf{T}}]^{\mathsf{T}}$. Intuitively, if $\mu$ is sufficiently large, the stability should be guaranteed. Indeed, the following proposition holds.

**Proposition 6** The closed-loop system with the controller designed through (4) with perturbed data $(Z_\Delta, X_\Delta, U_\Delta)$ is stable if

$$\mu > \left( -1 + \sqrt{1 + \frac{1}{2\|Z\|_2^2 \|M_\Delta\|_2}} \right)^{-1}$$

where $M_\Delta := G_\Delta P_\Delta G_\Delta^{\mathsf{T}}$ and $(P_\Delta, K_\Delta, G_\Delta)$ denotes the optimal solution to (4) under adversarial perturbation $\Delta$.

Proposition 6 establishes that a high signal-to-perturbation ratio guarantees stability under adversarial conditions. The regularization term $\|\Pi G\|$ in (4) serves to suppress the magnitude of the variables $G$ and $P$, which in turn reduces the spectral norm $\|M_\Delta\|_2$ and thereby enlarges the range of admissible values for the signal-to-perturbation ratio $\mu$. Furthermore, the proposition highlights the importance of selecting data with large $\sigma_{\min}(\Gamma)$ and $\|Z\|_2$, as such data enhance the robustness of the resulting controller against adversarial perturbations.

### B. Stability Enhancement by Robust Data-driven Control

We propose another defense method based on robust data-driven control [22], where the designed controller ensures stability under any realizable perturbation. While only robust stability is considered in [22], we additionally impose optimality on the quadratic performance for nominal data. We design the robust controller based on the following SDP:

$$\begin{aligned} \text{Find} \quad & L, S, \alpha, \beta \\ \text{s.t.} \quad & \mathrm{tr}(QX_\Delta L) + \mathrm{tr}(S) \le J_{\mathrm{LQR}}, \\ & F(L, S, D_\Delta) \succeq 0, \\ & \begin{bmatrix} X_\Delta L - \beta I_n & 0_{n,n} & 0_{n,m} & 0_{n,n} \\ 0_{n,n} & -X_\Delta L & -L^{\mathsf{T}} U_\Delta^{\mathsf{T}} & 0_{n,n} \\ 0_{m,n} & * & 0_{m,m} & U_\Delta L \\ 0_{n,n} & 0_{n,n} & * & X_\Delta L \end{bmatrix} - \alpha N_\Delta \succeq 0, \\ & \alpha \ge 0, \ \beta \ge 1 \end{aligned} \qquad (17)$$

with $J_{\mathrm{LQR}} > 0$,

$$N_\Delta := \mathrm{diag}\left( \bar{D}_\Delta^{\mathsf{T}} \begin{bmatrix} \epsilon^2(2n+m)TI_{2n+m} & 0 \\ 0 & -I_T \end{bmatrix} \bar{D}_\Delta, 0_{n,n} \right)$$

and $\bar{D}_\Delta := [I_{2n+m} \ [Z_\Delta^{\mathsf{T}} \ -X_\Delta^{\mathsf{T}} \ -U_\Delta^{\mathsf{T}}]^{\mathsf{T}}]^{\mathsf{T}}$. The resulting controller is given by $K = F_K(L, D_\Delta)$. The following proposition holds.

**Proposition 7** Consider the controller $K$ designed via (17). In the absence of perturbations (i.e., when $\Delta = 0$), the resulting controller achieves the nominal performance $\mathcal{J}(K) \le J_{\mathrm{LQR}}$. Furthermore, in the presence of perturbations, the closed-loop system remains stable for any $\Delta$ satisfying $\|\Delta\|_{\max} \le \epsilon$.

In contrast to the regularization-based defense strategy, the robust data-driven control approach offers a key advantage in that it provides a formal performance and stability guarantee. In addition, Proposition 7 guarantees that the proposed controller not only recovers optimal performance in the nominal

setting but also ensures robust stability under bounded perturbations, thereby bridging optimality and robustness within a unified design framework. Note that, in the relative constraint case with (6), $\epsilon$ in $N_\Delta$ is replaced with $\epsilon \max(\|D\|_{\max})$.

## VI. EXPERIMENTAL EVALUATION

### A. Experimental Setup

*Benchmarks:* To validate the generality and effectiveness of the proposed methods, we conduct numerical experiments on several benchmark LTI systems, including five of standard control models provided by the University of Michigan Control Tutorials for MATLAB and Simulink (CTMS) [24], [25], widely used in both education and research. Specifically, we consider the following systems: motor position control (MP), suspension system (SS), inverted pendulum (IP), aircraft pitch control (AP), and ball and beam (BB). These systems collectively exhibit a range of dynamic behaviors including instability, underactuation, and lightly damped oscillations. In addition, we include a linearized model of a triple tank system (TT) [26], [27], which represents a prototypical process control scenario involving interconnected fluid dynamics. The inclusion of TT complements the CTMS models and ensures that the proposed approach is applicable to a diverse array of LTI systems across multiple engineering domains. Although each system has a relatively small state-space dimension (fewer than five states), they serve as fundamental components in larger applications. Investigating their vulnerabilities therefore provides important insights into the reliability of data-driven control in broader real-world scenarios.

*Baseline Method:* To provide a baseline for comparison, we consider a *simple random attack strategy* that perturbs the data by sampling uniformly within the feasible perturbation set. Specifically, the perturbation matrix $\Delta$ is generated such that each entry is independently drawn from a uniform distribution over an interval that ensures $\|\Delta\|_{\max} \leq \epsilon$. This method does not exploit system structure or optimization but serves as a natural reference to assess the impact of more sophisticated, adversarially optimized attacks. By comparing against this baseline, we highlight the relative severity and effectiveness of targeted perturbations in degrading stability.

*Attack Impact Metric:* To quantitatively assess the effectiveness of adversarial perturbations, we adopt the *attack success rate (ASR)* as the primary metric. This ASR is defined as the ratio of the number of instances in which the controller designed from adversarially perturbed data results in an unstable closed-loop system to the total number of tested instances. This metric highlights how often an attack successfully destabilizes a previously stable system and serves as a direct indicator of the relative strength of different perturbation methods.

*Parameters:* All the systems are converted to discrete-time systems using an ideal sampler and a zero-order hold with the sampling period $T_s = 0.1s$. The time horizon is set to $T = 2(n + m)$. The input signal and the initial state are randomly and independently generated by the normal distribution, i.e., $u_t \sim \mathcal{N}(0, I_m)$ for $t = 0, \ldots, T - 1$ and $x_0 \sim \mathcal{N}(0, I_n)$. The state trajectory $x_t$ for $t = 1, \ldots, T$ is

TABLE I: Poles of open-loop system, closed-loop system with clean data, and closed-loop system with perturbed data.

| | Poles | | |
|---|---|---|---|
| Open-loop system | 0.9700 | 0.9600 | 0.9000 |
| Closed-loop (clean data) | 0.0012 | 0.0003 | 0.0001 |
| Closed-loop (perturbed data) | 1.0016 | 0.4723 | 0.1550 |

generated by the dynamics. The LQR weight matrices are set to $Q = I_n$ and $R = 0.01I_m$. The step size at each iteration of I-DGSM is set to $\alpha_{\text{step}} = \epsilon \|D\|_{\max}$. The termination condition of I-DGSM is met when either $\rho(A + B\mathcal{K}(D_k)) \geq 1$ or the iteration count exceeds 50, i.e., $k > 50$. The relative perturbation constraint (6) is adopted in all instances.

*Computational Environment:* All numerical experiments were performed on a workstation with 3.10 GHz and 12 cores CPU and 64 GB RAM. Computations were implemented in MATLAB® (version R2024b) [46]. The source code used in the numerical experiments is available online[1].

### B. Visualization

First of all, we illustrate the adversarial attack. Fig. 3 presents a visualization of the proposed adversarial attack applied to TT using I-DGSM with $\epsilon = 0.005$ where the regularization parameter is set to $\gamma = 0.0001$. In the upper half, the clean data, the adversarial perturbation, and the perturbed signal of the first input and state signals are depicted. In the lower left, the eigenvalues of the resulting closed-loop system with the clean data are depicted, while in the lower right, those with the perturbed data are shown. The illustration highlights that while the system can be stabilized using clean data, resulting in eigenvalues significantly distant from the unit circle, a small yet sophisticated perturbation can render the system unstable. The poles of the open-loop system, those of the closed-loop system with the clean data, and those of the closed-loop system with the perturbed data are summarized in TABLE I. As observed, the controller designed using clean data yields poles near the origin, indicating strong stability, whereas the adversarially perturbed data leads to a pole outside the unit circle, causing instability. Notably, the perturbed signals appear visually similar to the clean ones, suggesting that detecting such malicious modifications through inspection alone would be challenging. This result underscores the vulnerability of direct data-driven control methods to adversarial data manipulation.

### C. Attack Impact Evaluation

Next, we evaluate the attack impact of the three perturbation generation algorithms. Fig. 4 shows ASR over 20 trials as a function of the perturbation size $\epsilon$ for six benchmark systems, using the baseline, DGSM, and I-DGSM algorithms, where the regularization parameter is set to $\gamma = 1$. The dotted, dashed, and solid lines represent the baseline, DGSM, and I-DGSM, respectively. The horizontal axis is shown on a

---

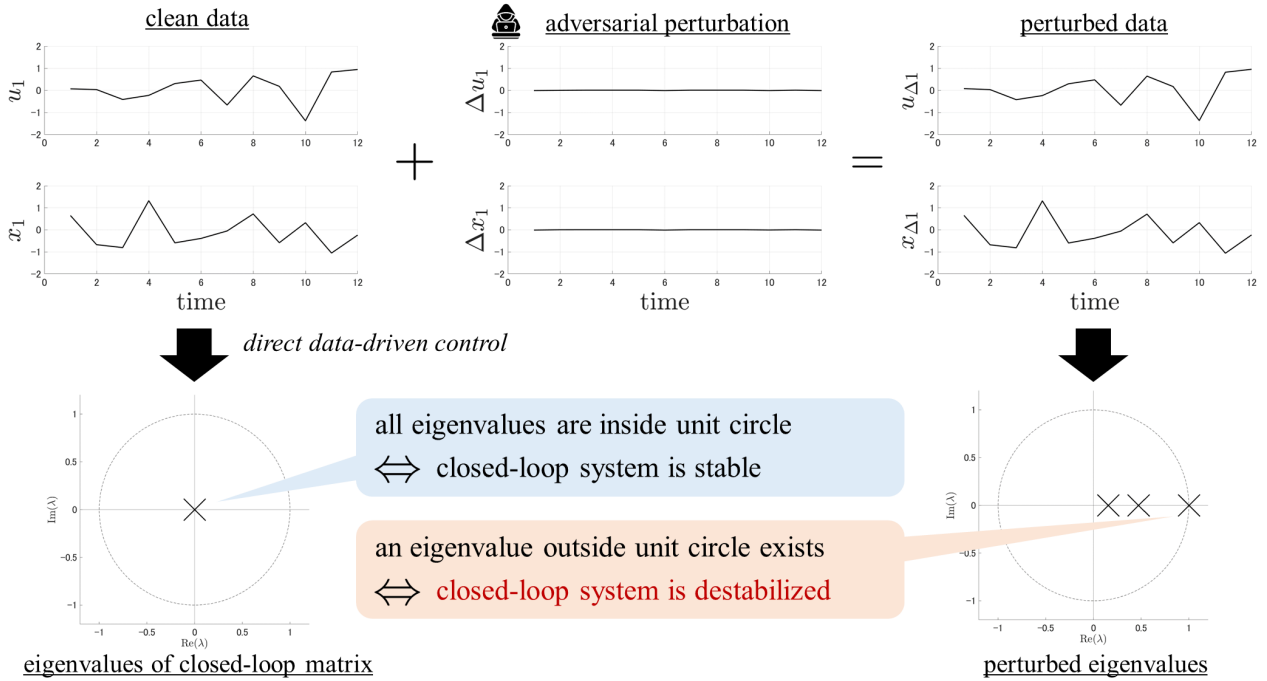[1]https://github.com/HampeiSasahara/adv-des-attack-dddc

Fig. 3: Visualization of the adversarial attack. In the upper half, the clean data, the adversarial perturbation, and the perturbed signal of the first input and state signals are depicted. In the lower half, the eigenvalues of the closed-loop matrix where the controller is designed through the offline direct data-driven control with the clean data and those with the perturbed data are depicted. With the clean data, all eigenvalues are inside the unit circle, indicating stability of the resulting closed-loop system. Despite the addition of an adversarial perturbation to the original signal, the perturbed signals appear almost identical. However, when using the perturbed data, an eigenvalue outside the unit circle is observed, indicating destabilization of the closed-loop system due to the adversarial attack. Note that the second and third input and state signals are also perturbed but its illustration is omitted for clarity.

logarithmic scale to highlight differences across a wide range of perturbation magnitudes.

This result reveals several consistent patterns. First, across all benchmark systems, I-DGSM consistently induces instability with the smallest perturbation sizes, demonstrating its effectiveness in identifying the most vulnerable directions. DGSM requires slightly larger perturbations to cause instability but still outperforms the baseline method in most cases. The baseline approach, in contrast, requires significantly larger perturbations to achieve similar effects. These trends are in line with expectations, as I-DGSM incorporates optimization that explicitly targets the most destabilizing perturbations.

Furthermore, the performance gap between I-DGSM and the baseline is substantial. In MP, IP, BB, and TT, I-DGSM achieves a comparable ASR using perturbations that are approximately $10^{-1}$ times smaller. This gap increases to around $10^{-2}$ in SS and AP. These results underscore the efficiency and precision of I-DGSM in identifying minimal but impactful perturbations for destabilizing the system.

### D. Computational Cost Evaluation

We here evaluate the effectiveness of the proposed gradient computation method. TABLE II compares the computation times in seconds (mean $\pm$ standard deviation over 10 trials) required to evaluate the gradient $\nabla_D \rho(A + B\mathcal{K}(D))$ under

TABLE II: Computation time (in seconds) for evaluating gradient $\nabla_D \rho(A + B\mathcal{K}(D))$.

|     | NumDiff | ImpDiffNum | ImpDiffAnal |
| --- | --- | --- | --- |
| MP | 41.874$\pm$2.234 | 0.613$\pm$0.036 | 0.380$\pm$0.025 |
| SS | 99.730$\pm$6.668 | 1.373$\pm$0.051 | 0.470$\pm$0.044 |
| IP | 83.401$\pm$1.531 | 1.125$\pm$0.062 | 0.470$\pm$0.026 |
| AP | 54.695$\pm$1.868 | 0.761$\pm$0.096 | 0.479$\pm$0.035 |
| BB | 81.427$\pm$1.435 | 1.135$\pm$0.061 | 0.473$\pm$0.028 |
| TT | 82.322$\pm$5.130 | 1.087$\pm$0.106 | 0.411$\pm$0.038 |

different differentiation methods: NumDiff refers to numerical differentiation where the entire gradient is computed using the central difference method [38, Chap. 4]. ImpDiffNum refers to implicit differentiation with numerical derivatives where the overall gradient is computed via the implicit differentiation as described in Algorithm 2 but the required derivatives in (8) and (12) are evaluated numerically using central difference method. ImpDiffAnal refers to implicit differentiation with analytical derivatives where both the gradient and the necessary intermediate derivatives are computed using the closed-form expressions provided in Section IV-C.

A substantial improvement is observed when moving from NumDiff to ImpDiffNum. Across all six instances, the computation time is reduced by factors ranging from approximately 68$\times$ (MP: 41.87 s to 0.61 s) to 75$\times$ (TT: 82.32 s to 1.09 s),
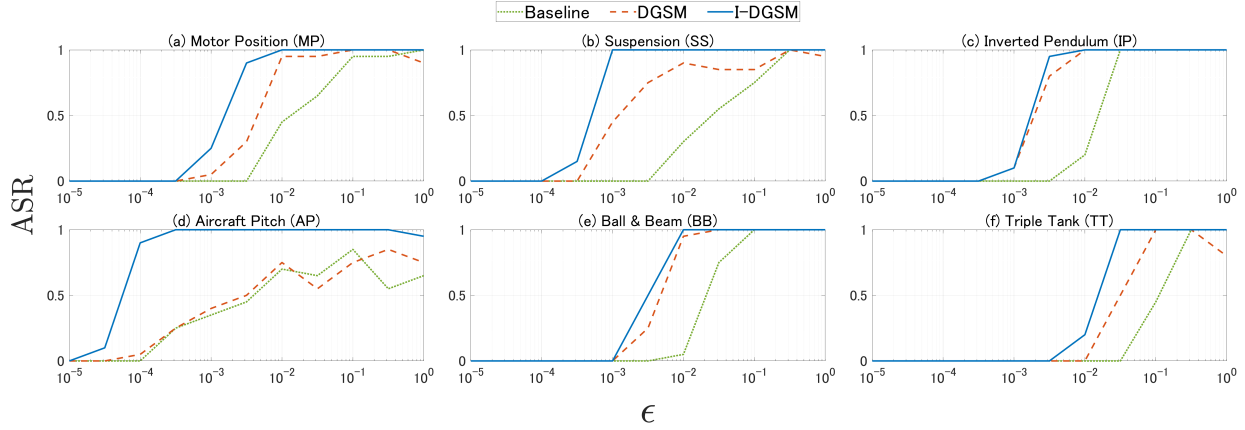
Fig. 4: ASR over 20 trials as a function of the perturbation size $\epsilon$ (shown on the horizontal axis in log scale) for the three perturbation generation methods (baseline, DGSM, and I-DGSM) across six benchmark systems. Subfigure titles indicate the specific system.

with other cases such as IP, SP, AP, and BB also showing speedups of $74\times$, $73\times$, $72\times$, and $72\times$, respectively. This dramatic reduction confirms that even when numerical differentiation is still used to compute intermediate quantities, implicit differentiation offers a significant computational advantage by avoiding redundant gradient evaluations. An additional improvement is achieved by switching from numerical to analytical evaluation of the derivatives, i.e., from ImpDiffNum to ImpDiffAnal. The resulting speedups range from about $1.6\times$ (MP: 0.61 s to 0.38 s) to $2.9\times$ (SS: 1.37 s to 0.47 s), with similar gains for other systems (IP: $2.4\times$, AP: $1.6\times$, BB: $2.4\times$, TT: $2.6\times$). This further highlights the importance of analytical formulations in maximizing the efficiency of gradient computations, particularly in settings where repeated evaluations for large-scale systems are required.

### E. Effectiveness of Defense Strategies

Subsequently, we evaluate the effectiveness of the two proposed defense strategies.

*1) Regularization:* Since the appropriate range of regularization parameters can vary depending on the system, we begin with a detailed analysis of the IP system as a representative example. The top panel of Fig. 5 shows ASR over 20 trials as a function of the regularization parameter $\gamma$ for the IP system. We can observe a clear trend that the ASR decreases as $\gamma$ increases. Further, to assess the impact of regularization on control performance, we also evaluate the averaged relative control performance (RCP) $\mathcal{J}(K_{\mathrm{reg}})/\mathcal{J}(K_{\mathrm{LQR}})$, where $K_{\mathrm{reg}}$ and $K_{\mathrm{LQR}}$ denote the controller designed by (5) and the ideal LQR optimal controller, respectively. Since the SDP (5) yields the exact solution under clean data, we add random noise (uniformly distributed in $[-\epsilon, \epsilon]$), with which the designed controller maintains the stability, to simulate perturbed data. The bottom panel of Fig. 5 illustrates that the control performance improves as the regularization becomes stronger under noisy data. Specifically, RCP decreases from 6.22 at $\gamma = 10^{-4}$ to 1.00 at $\gamma = 10^2$. TABLE III summarizes these trends for the other systems, where we define $\epsilon_{\mathrm{e}} := \log_{10} \epsilon$
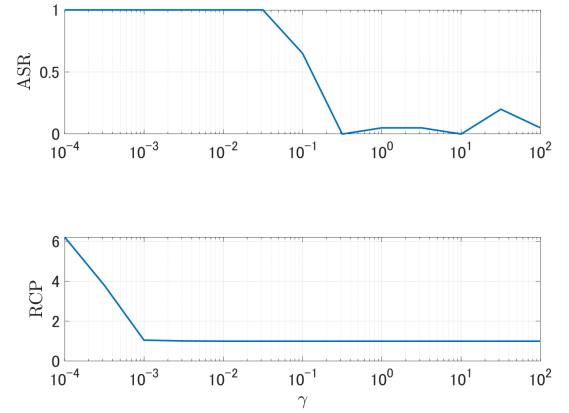


Fig. 5: ASR and mean of RCP $\mathcal{J}(K_{\mathrm{reg}})/\mathcal{J}(K_{\mathrm{LQR}})$ over 20 trials as a function of the regularization parameter $\gamma$ (shown on the horizontal axis in log scale) for IP with the perturbation size $\epsilon = 0.001$.

TABLE III: Effectiveness of regularization-based defense strategy.

| | $\epsilon_{\mathrm{e}}$ | Varied $\gamma_{\mathrm{e}}$ | Varied ASR | Varied RCP |
|---|---|---|---|---|
| MP | $-3$ | $-2 \to +0$ | $1.00 \to 0.30$ | $1.00 \to 1.00$ |
| SS | $-4$ | $-4 \to -2$ | $1.00 \to 0.00$ | $1.00 \to 1.00$ |
| IP | $-3$ | $-2 \to +0$ | $1.00 \to 0.05$ | $1.00 \to 1.00$ |
| AP | $-5$ | $-5 \to -3$ | $1.00 \to 0.00$ | $1.00 \to 1.00$ |
| BB | $-3$ | $-3 \to -1$ | $1.00 \to 0.05$ | $1.03 \to 1.00$ |
| TT | $-2$ | $-5 \to -3$ | $1.00 \to 0.30$ | $1.06 \to 1.00$ |

and $\gamma_{\mathrm{e}} := \log_{10} \gamma$. In all cases, the ASR decreases with increasing $\gamma$, while the RCP remains largely stable or improves, confirming the robustness benefits of regularization. These results collectively demonstrate the effectiveness of the proposed regularization-based defense method.

*2) Robust Data-driven Control:* We next evaluate the effectiveness of the robust data-driven control strategy. Since

TABLE IV: RCP $\mathcal{J}(K_{\mathrm{rob}})/\mathcal{J}(K_{\mathrm{LQR}})$ with robust data-driven control for various $\epsilon_{\mathrm{e}}$.

| | $\mathcal{J}(K_{\mathrm{rob}})/\mathcal{J}(K_{\mathrm{LQR}})$ for | | | | |
| | $\epsilon_{\mathrm{e}} = -6$ | $-5$ | $-4$ | $-3$ | $-2$ |
|---|---|---|---|---|---|
| MP | 1.05 | 1.42 | N/A | N/A | N/A |
| SS | 1.01 | 1.24 | 5.62 | N/A | N/A |
| IP | 1.00 | 1.01 | 1.38 | N/A | N/A |
| AP | 1.20 | 3.21 | N/A | N/A | N/A |
| BB | 1.00 | 1.01 | 1.17 | 2.86 | N/A |
| TT | 1.00 | 1.00 | 1.01 | N/A | N/A |

this method provides formal guarantees of closed-loop stability under admissible perturbations, we focus here only on assessing the control performance degradation. Specifically, we compute the RCP, defined as $\mathcal{J}(K_{\mathrm{rob}})/\mathcal{J}(K_{\mathrm{LQR}})$, where $K_{\mathrm{rob}}$ is the controller obtained from the robust data-driven formulation (17) while minimizing $J_{\mathrm{LQR}}$. TABLE IV reports the RCP values for different levels of perturbation, indexed by $\epsilon_{\mathrm{e}}$. Entries marked as N/A indicate that the SDP (17) is infeasible.

The results show that RCP remains close to 1 when $\epsilon$ is on the order of $10^{-6}$, indicating minimal performance degradation. However, as $\epsilon$ increases, the RCP grows significantly, and the problem becomes infeasible for $\epsilon$ around $10^{-3}$ to $10^{-2}$, which is a range where the attack achieves high ASR in most systems as shown in Fig. 4. In comparison with TABLE III, these results suggest that it is less effective in preserving control performance than the regularization-based defense strategy while the robust data-driven control strategy provides formal stability guarantees.

### F. Transferability Evaluation

Finally, we investigate transferability across data, where the adversary does not have access to the dataset $D$ used for controller design. Instead, the adversary generates a hypothetical dataset $D_{\mathrm{hyp}}$ based on synthetic input-state trajectories. Specifically, the hypothetical input is sampled as $u_{\mathrm{hyp},t} \sim \mathcal{N}(0, I_m)$, with an initial state $x_{\mathrm{hyp},0} \sim \mathcal{N}(0, I_n)$, and the resulting state trajectory $x_{\mathrm{hyp},t}$ is generated for $t = 1, \ldots, T$ by simulating the system dynamics.

Fig. 6 illustrates ASR over 20 trials as a function of the perturbation magnitude $\epsilon$ for three scenarios (baseline, full knowledge, and partial knowledge) across six benchmark systems with $\gamma = 1$. From these graphs, we observe that the partial knowledge attacks typically require perturbations with magnitudes about $10^1$ to $10^2$ times larger than those in the full knowledge case to achieve comparable ASR. Nevertheless, the partial knowledge attacks are more effective than the baseline in most cases. These results indicate that data transferability is moderate under this threat model and suggest that maintaining the confidentiality of the training data can serve as a practical defense against adversarial perturbations in control systems.

### VII. CONCLUSION

In this study, we have revealed the vulnerability of direct data-driven control to adversarial attacks, focusing on destabilizing closed-loop systems in the LQR setting. We have proposed effective perturbation generation algorithms, namely DGSM and I-DGSM, tailored for control systems. To enable practical attack evaluation, we have developed an efficient gradient computation method based on implicit differentiation through the KKT conditions of the underlying semidefinite program. For defense, we have introduced two proactive strategies: a regularization-based approach that enhances stability against adversarial perturbations, and a robust data-driven control method guaranteeing closed-loop stability under bounded data perturbations. Extensive numerical experiments have demonstrated the severity of the attacks, the computational efficiency of our methods, and the effectiveness of the proposed defenses. Furthermore, transferability analysis under partial knowledge scenarios has highlighted the importance of training data confidentiality as a practical defense.

Future work includes extending the proposed attack and defense frameworks to nonlinear and time-varying systems, as well as developing real-time detection and mitigation mechanisms suitable for online operation. Investigating robustness under more realistic threat models with limited attacker capabilities remains important directions. Additionally, resolving the issue on storing and manipulating large matrices is also included in the future work. Finally, leveraging the efficient gradient computation techniques presented here, we aim to address safety and security challenges in large-scale and networked control systems, ultimately contributing to the deployment of resilient data-driven controllers in complex, safety-critical environments.

### APPENDIX A
### BASIC RULES OF MATRIX CALCULUS

Following the basic matrix calculus [44], we define the derivative of a matrix-valued function $F(X) = [f_1 \; \cdots \; f_m] \in \mathbb{R}^{n \times m}$ with respect to the matrix variable $X = [x_1 \; \cdots \; x_q] \in \mathbb{R}^{p \times q}$ by

$$\frac{dF}{dX} := \frac{d\mathrm{vec}(F)}{d\mathrm{vec}(X)} = \begin{bmatrix} \dfrac{df_1}{dx_1} & \cdots & \dfrac{df_1}{dx_q} \\ \vdots & & \vdots \\ \dfrac{df_m}{dx_1} & \cdots & \dfrac{df_m}{dx_q} \end{bmatrix} \in \mathbb{R}^{nm \times pq},$$

(18)

where $df_i/dx_j \in \mathbb{R}^{n \times p}$ denotes the corresponding Jacobian matrix. We define the partial derivative $\partial F/\partial X$ in a similar manner. We define the gradient of a scalar function $f(X)$ with respect to the matrix variable $X \in \mathbb{R}^{p \times q}$ by $\nabla_X f = \mathrm{vec}^{-1}((df/dX)^{\mathsf{T}}) \in \mathbb{R}^{p \times q}$. Based on the definition (18), we have the chain rule for $H(X) = G(F(X))$ as

$$dH/dX = dG/dY|_{Y=F(X)} \, dF/dX$$

and the product rule for $H(X) = F(X)G(X)$ with $F(X) \in \mathbb{R}^{n \times m}, G(X) \in \mathbb{R}^{m \times \ell}$ as

$$dH/dX = (G^{\mathsf{T}} \otimes I_n)dF/dX + (I_\ell \otimes F)dG/dX.$$

TABLE V exhibits useful formulae given $F(X) \in \mathbb{R}^{n \times m}$ with $X \in \mathbb{R}^{p \times q}$ [44, Chap. 9].
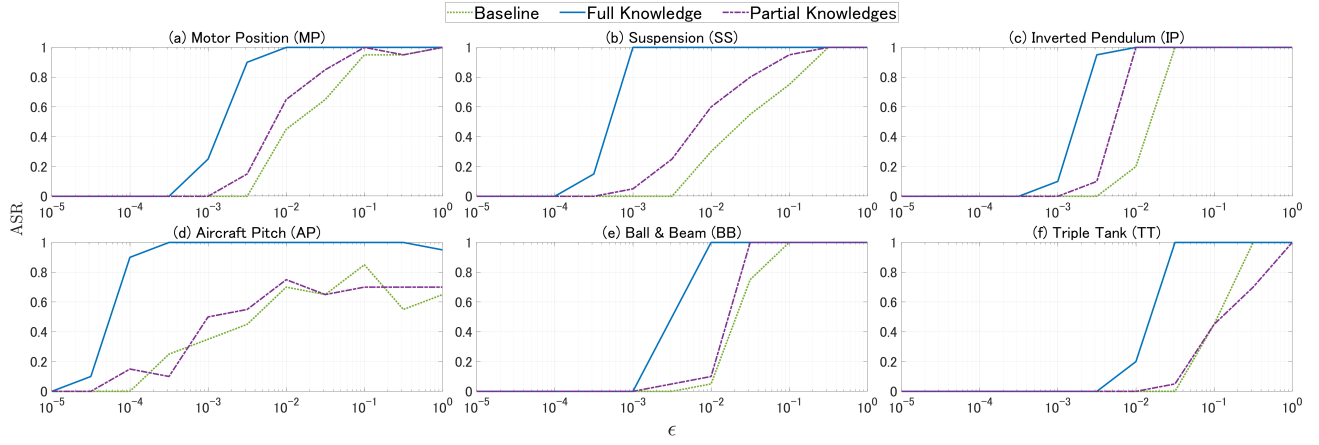
Fig. 6: ASR over 20 trials as a function of the perturbation size $\epsilon$ (shown on the horizontal axis in log scale) for the three cases (baseline, full knowledge, and partial knowledge) across six benchmark systems. Subfigure titles indicate the specific system.

TABLE V: Matrix derivative formulae.

| $F$ | $dF/dX$ |
|---|---|
| $X^{\mathsf{T}}$ | $C_{p,q}$ |
| $X^{-1}$ | $-X^{-\mathsf{T}} \otimes X^{-1}$ |
| $AXB$ | $B^{\mathsf{T}} \otimes A$ |
| $\mathrm{vec}(X)^{\mathsf{T}}$ | $I_{pq}$ |
| $\mathrm{tr}(AXB)$ | $(\mathrm{vec}(BA)^{\mathsf{T}})^{\mathsf{T}}$ |
| $\mathrm{tr}(X^{\mathsf{T}}AX)$ | $(\mathrm{vec}(A+A^{\mathsf{T}})X)^{\mathsf{T}}$ |
| $A \otimes X$ with $A \in \mathbb{R}^{r \times s}$ | $(I_s \otimes C_{q,r} \otimes I_p)(\mathrm{vec}(A) \otimes I_{pq})$ |
| $X \otimes A$ with $A \in \mathbb{R}^{r \times s}$ | $(I_q \otimes C_{s,p} \otimes I_r)(I_{pq} \otimes \mathrm{vec}(A))$ |

## APPENDIX B
## PROOFS

*Proof of Proposition 1:* Let us focus on the fourth block row of $H$, namely

$$[0_{(3n+m)^2,nT}\ 0_{(3n+m)^2,m^2}\ 0_{(3n+m)^2,1}\ I_{(3n+m)^2}]H = [\partial G_4/\partial L\ \ \partial G_4/\partial S\ \ \partial G_4/\partial \Lambda]. \quad (19)$$

Recall that $G_4 = \Lambda - \Lambda^{\mathsf{T}}$. Since the first diagonal component of $G_4$ is always zero, the first row of (19) is also zero given any variables. Similarly, since the other diagonal components are always zero, and the corresponding rows are also zero. Because the dimension of $G_4$ is $(3n+m) \times (3n+m)$, $H$ contains $3n+m$ zero row vectors at least. Thus, the rank of $H$ is less than $nT + m^2 + 1 + (3n+m)^2 - (3n+m) < nT + m^2 + (3n+m)^2$, which the number of columns. This leads to the claim. $\square$

*Proof of Proposition 2:* The sufficiency is obvious. We here show the necessity. Consider the Cholesky decomposition [47, Fact 10.10.42] $S_i = L_i L_i^{\mathsf{T}}$ with lower triangular matrices $L_i$ for $i = 1, 2$. Because $S_1 \succeq 0$, we have $L_2^{\mathsf{T}} S_1 L_2 \succeq 0$. Thus all eigenvalues of $L_2^{\mathsf{T}} S_1 L_2$ are nonnegative. Now $\mathrm{tr}(L_2^{\mathsf{T}} S_1 L_2) = \mathrm{tr}(S_1 L_2 L_2^{\mathsf{T}}) = \mathrm{tr}(S_1 S_2) = 0$. Because trace is equal to the sum of all eigenvalues, all eigenvalues of $L_2^{\mathsf{T}} S_1 L_2$ are zero, which implies $L_2^{\mathsf{T}} S_1 L_2$ is nilpotent. Since $L_2^{\mathsf{T}} S_1 L_2$ is symmetric, we have $L_2^{\mathsf{T}} S_1 L_2 = 0$. Then $L_2^{\mathsf{T}} S_1 L_2 = L_2^{\mathsf{T}} L_1 L_1^{\mathsf{T}} L_2 = (L_1^{\mathsf{T}} L_2)^{\mathsf{T}} L_1^{\mathsf{T}} L_2 = 0$, which means $L_1^{\mathsf{T}} L_2 = 0$. Therefore, $S_1 S_2 = L_1 L_1^{\mathsf{T}} L_2 L_2^{\mathsf{T}} = 0$. $\square$

*Proof of Proposition 3:* Consider $d\rho/dK$. Note that $d|\lambda|/dK = \mathrm{Re}(\lambda\, \overline{d\lambda/dK})/|\lambda|$. Since $(A + BK)v = v\lambda$ and $w^{\mathsf{T}}(A + BK) = \lambda w$, from the chain rule, we have

$$(v^{\mathsf{T}} \otimes I_n)\underbrace{d(A + BK)/dK}_{=I_n \otimes B} + A\,dv/dK = \lambda\,dv/dK + v\,d\lambda/dK.$$

Multiplying $w^{\mathsf{T}}$ from left yields

$$w^{\mathsf{T}}(v^{\mathsf{T}} \otimes B) + \underbrace{w^{\mathsf{T}}A}_{=\lambda w^{\mathsf{T}}}\,dv/dK = \lambda w^{\mathsf{T}}\,dv/dK + w^{\mathsf{T}}v\,d\lambda/dK,$$

and hence $w^{\mathsf{T}}(v^{\mathsf{T}} \otimes I_n)d(A + BK)/dK = w^{\mathsf{T}}v\,d\lambda/dK$. Since $w^{\mathsf{T}}v \neq 0$, we obtain $d\lambda/dK = w^{\mathsf{T}}(v^{\mathsf{T}} \otimes B)/(w^{\mathsf{T}}v)$.

Next, from the product rule, the chain rule, and TABLE V, we have

$$\partial F_K/\partial L = ((XL)^{-\mathsf{T}} \otimes I_m)(I_n \otimes U) - ((XL)^{-\mathsf{T}} \otimes F_K)(I_n \otimes X),$$
$$\partial F_K/\partial D = ((XL)^{-\mathsf{T}} \otimes I_m)(L^{\mathsf{T}} \otimes E_U) - ((XL)^{-\mathsf{T}} \otimes F_K)(L^{\mathsf{T}} \otimes E_X),$$

which lead to the analytic forms. $\square$

*Proof of Proposition 4:* Since $\|\Pi L\|_{\mathrm{F}}^2 = \mathrm{tr}(L^{\mathsf{T}}\Pi^{\mathsf{T}}\Pi L) = \mathrm{tr}(L^{\mathsf{T}}\Pi L)$, from TABLE V we have

$$\partial J/\partial L = (\mathrm{vec}(X^{\mathsf{T}}Q^{\mathsf{T}}))^{\mathsf{T}} + 2\gamma(\mathrm{vec}(\Pi L))^{\mathsf{T}},$$
$$\partial J/\partial S = (\mathrm{vec}(I_m))^{\mathsf{T}},$$

which lead to (13). Furthermore, because $F = \bar{E}_1 F_1 \bar{E}_1^{\mathsf{T}} + \bar{E}_2 F_2 \bar{E}_2^{\mathsf{T}}$, we have (14). Finally, considering the decomposition

$$F_1 = \begin{bmatrix} I_m \\ 0_{n,m} \end{bmatrix} S \begin{bmatrix} I_m\ 0_{m,n} \end{bmatrix} + \begin{bmatrix} I_m \\ 0_{n,m} \end{bmatrix} VUL \begin{bmatrix} 0_{n,m}\ I_n \end{bmatrix} + \begin{bmatrix} 0_{m,n} \\ I_n \end{bmatrix} (VUL)^{\mathsf{T}} \begin{bmatrix} I_m\ 0_{m,n} \end{bmatrix} + \begin{bmatrix} 0_{m,n} \\ I_n \end{bmatrix} XL \begin{bmatrix} 0_{n,m}\ I_n \end{bmatrix} \quad (20)$$

and that for $F_2$ in a similar manner, we obtain (15). $\square$

*Proof of Proposition 5:* Using the decompositions of $F_1$, $F_2$, $\partial F_1/\partial L$, and $\partial F_2/\partial L$ as in (20), together with the formulae in TABLE V, we obtain the desired derivatives. For the

derivative $d\Pi/dD$, we apply the identity $\Gamma^\dagger = \Gamma^\mathsf{T}(\Gamma\Gamma^\mathsf{T})^{-1}$ since $\Gamma$ is a full-row rank matrix. □

*Proof of Proposition 6:* Since $Z = [B\ A]\Gamma$ and $\Gamma$ is full-row rank, we have $[B\ A] = Z\Gamma^\dagger$ and $Z = [B\ A](\Gamma_\Delta - \Delta_\Gamma)$. Thus, from (1),

$$
\begin{aligned}
ZG_\Delta &= [B\ A](\Gamma_\Delta - \Delta_\Gamma)G_\Delta \\
&= [B\ A]([K_\Delta^\mathsf{T}\ I_n]^\mathsf{T} - \Delta_\Gamma G_\Delta) \\
&= A + BK_\Delta - Z\Gamma^\dagger\Delta_\Gamma G_\Delta,
\end{aligned}
$$

which leads to the closed-loop matrix expression $A + BK_\Delta = Z(I + \Gamma^\dagger\Delta_\Gamma)G_\Delta$. Hence, the stability condition can be characterized by the Lyapunov inequality $\mathrm{Lyap}(P) \prec 0$ for some $P \succ 0$ where

$$
\mathrm{Lyap}(P) := Z(I + \Gamma^\dagger\Delta_\Gamma)G_\Delta P G_\Delta^\mathsf{T}(I + \Gamma^\dagger\Delta_\Gamma)^\mathsf{T}Z^\mathsf{T} - P.
$$

Let us take $P_\Delta$ as a candidate that satisfies the Lyapunov inequality. Note that, from the constraint in (4), $P_\Delta$ satisfies $Z_\Delta M_\Delta Z_\Delta^\mathsf{T} - P_\Delta + I_n \preceq 0$. Now we define $\Psi := \mathrm{Lyap}(P_\Delta) - Z_\Delta M_\Delta Z_\Delta^\mathsf{T} + P_\Delta$. Using basic algebra, we have

$$
\begin{aligned}
\Psi = {}&- ZM_\Delta\Delta_Z^\mathsf{T} - \Delta_Z M_\Delta Z^\mathsf{T} - \Delta_Z M_\Delta\Delta_Z^\mathsf{T} + Z\Gamma^\dagger\Delta_\Gamma M_\Delta Z^\mathsf{T} \\
&+ ZM_\Delta(\Gamma^\dagger\Delta_\Gamma)^\mathsf{T}Z^\mathsf{T} + Z\Gamma^\dagger\Delta_\Gamma M_\Delta(\Gamma^\dagger\Delta_\Gamma)^\mathsf{T}Z^\mathsf{T}.
\end{aligned}
$$

Because every operator norm is submultiplicative and $\|\Gamma^\dagger\|_2 = \sigma_{\min}(\Gamma)^{-1}$, we have

$$
\begin{aligned}
\|\Psi\|_2 \leq {}&\|Z\|_2^2\|M_\Delta\|_2(\|\Delta_Z\|_2^2/\|Z\|_2^2 + 2\|\Delta_Z\|_2/\|Z\|_2 \\
&+ \|\Delta_\Gamma\|_2^2/\sigma_{\min}(\Gamma)^2 + 2\|\Delta_\Gamma\|_2/\sigma_{\min}(\Gamma)) \\
\leq {}&2\|Z\|_2^2\|M_\Delta\|_2(\mu^{-2} + 2\mu^{-1}).
\end{aligned}
$$

Therefore, if $\mu$ satisfies the inequality, we have $\|\Psi\|_2 \leq 1$, which is equivalent to $\Psi \prec I$. This implies that $\mathrm{Lyap}(P_\Delta) \prec Z_\Delta M_\Delta Z_\Delta^\mathsf{T} - P_\Delta + I \prec 0$, which leads to the claim. □

*Proof of Proposition 7:* We first consider the perturbation-free case, i.e., $\Delta = 0$. In this case, it follows that $D_\Delta = D$ and $Z_\Delta = AX_\Delta + BU_\Delta$. Define $K = F_K(L, D_\Delta)$, $P = X_\Delta L$, and $G = LP^{-1}$, where $P$ is nonsingular due to the constraint $F_2 \succeq 0$. By applying the Schur complement to the condition $F(L, S, D_\Delta) \succeq 0$, it can be verified that the triplet $(K, P, G)$ satisfies the constraints in (3). Therefore, the nominal performance is guaranteed by $\mathcal{J}(K) = \mathrm{trace}(QX_\Delta L) + \mathrm{trace}(S) \leq J_{\mathrm{LQR}}$.

Next, consider the case where $\|\Delta\|_{\max} \leq \epsilon$. We have $\|\Delta\|_2 \leq \|\Delta\|_F \leq \sqrt{(2n+m)T}\|\Delta\|_{\max} \leq \sqrt{(2n+m)T}\epsilon$, which implies $\Delta\Delta^\mathsf{T} \preceq \epsilon^2(2n+m)TI_{2n+m}$. This is equivalent to the quadratic matrix inequality $[I\ \Delta]\mathrm{diag}(\epsilon^2(2n+m)TI_{2n+m}, -I_T)[I\ \Delta]^\mathsf{T} \preceq 0$. According to Theorem 3.8 in [22], the third, fourth, and fifth constraints in (17) constitute a sufficient condition for the stability of the closed-loop system.

## REFERENCES

[1] W. Schwarting, J. Alonso-Mora, and D. Rus, "Planning and decision-making for autonomous vehicles," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 1, no. 1, pp. 187–210, 2018.

[2] J. Bruna, C. Szegedy, I. Sutskever, I. Goodfellow, W. Zaremba, R. Fergus, and D. Erhan, "Intriguing properties of neural networks," in *Proc. International Conference on Learning Representations (ICLR)*, 2014.

[3] I. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *Proc. International Conference on Learning Representations (ICLR)*, 2015.

[4] N. Carlini, P. Mishra, T. Vaidya, Y. Zhang, M. Sherr, C. Shields, D. Wagner, and W. Zhou, "Hidden voice commands," in *25th USENIX security symposium (USENIX security 16)*, 2016, pp. 513–530.

[5] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, "Robust physical-world attacks on deep learning visual classification," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1625–1634.

[6] J. Tian, B. Wang, Z. Wang, K. Cao, J. Li, and M. Ozay, "Joint adversarial example and false data injection attacks for state estimation in power systems," *IEEE Trans. Cybern.*, vol. 52, no. 12, pp. 13 699–13 713, 2021.

[7] Y. Wang, E. Sarkar, W. Li, M. Maniatakos, and S. E. Jabari, "Stop-and-go: Exploring backdoor attacks on deep reinforcement learning-based traffic congestion control systems," *IEEE Trans. Inf. Forensics Security*, vol. 16, pp. 4772–4787, 2021.

[8] S. A. Deka, D. M. Stipanović, and C. J. Tomlin, "Dynamically computing adversarial perturbations for recurrent neural networks," *IEEE Trans. Control Syst. Technol.*, vol. 30, no. 6, pp. 2615–2629, 2022.

[9] R. Song, M. O. Ozmen, H. Kim, R. Muller, Z. B. Celik, and A. Bianchi, "Discovering adversarial driving maneuvers against autonomous vehicles," in *Proc. 32nd USENIX Security Symposium (USENIX Security 23)*, 2023, pp. 2957–2974.

[10] E. O'Dwyer, E. C. Kerrigan, P. Falugi, M. Zagorowska, and N. Shah, "Data-driven predictive control with improved performance using segmented trajectories," *IEEE Trans. Control Syst. Technol.*, vol. 31, no. 3, pp. 1355–1365, 2022.

[11] F. Dörfler, "Data-driven control: Part one of two: A special issue sampling from a vast and dynamic landscape," *IEEE Control Systems Magazine*, vol. 43, no. 5, pp. 24–27, 2023.

[12] L. Schmitt, J. Beerwerth, M. Bahr, and D. Abel, "Data-driven predictive control with online adaption: Application to a fuel cell system," *IEEE Trans. Control Syst. Technol.*, vol. 32, no. 1, pp. 61–72, 2023.

[13] M. C. Campi, A. Lecchini, and S. M. Savaresi, "Virtual reference feedback tuning: A direct method for the design of feedback controllers," *Automatica*, vol. 38, no. 8, pp. 1337–1346, 2002.

[14] O. Kaneko, "Data-driven controller tuning: FRIT approach," *IFAC Proceedings Volumes*, vol. 46, no. 11, pp. 326–336, 2013.

[15] F. L. Lewis and D. Liu, *Reinforcement Learning and Approximate Dynamic Programming for Feedback Control*. John Wiley & Sons, 2013.

[16] Y. Yu, R. Zhao, S. Chinchali, and U. Topcu, "Poisoning attacks against data-driven predictive control," in *Proc. American Control Conference (ACC)*, 2023, pp. 545–550.

[17] T. Ikezaki, O. Kaneko, K. Sawada, and J. Fujita, "Poisoning attack on VIMT and its adverse effect," *Artificial Life and Robotics*, vol. 29, no. 1, pp. 168–176, 2024.

[18] F. Fotiadis, A. Kanellopoulos, K. G. Vamvoudakis, and U. Topcu, "Deception against data-driven linear-quadratic control," *arXiv preprint*, 2025, [Online]. Available: https://arxiv.org/pdf/2506.11373.

[19] C. De Persis and P. Tesi, "Formulas for data-driven control: Stabilization, optimality, and robustness," *IEEE Trans. Autom. Control*, vol. 65, no. 3, pp. 909–924, 2019.

[20] F. Dörfler, P. Tesi, and C. De Persis, "On the role of regularization in direct data-driven LQR control," in *Proc. IEEE 61st Conference on Decision and Control (CDC)*, 2022, pp. 1091–1098.

[21] F. Dörfler, J. Coulson, and I. Markovsky, "Bridging direct & indirect data-driven control formulations via regularizations and relaxations," *IEEE Trans. Autom. Control*, vol. 68, no. 2, pp. 883–897, 2023.

[22] T. Kaminaga and H. Sasahara, "Data informativity under data perturbation," *arXiv preprint*, 2025, [Online]. Available: https://arxiv.org/pdf/2505.01641.

[23] ——, "Data informativity for quadratic stabilization under data perturbation," in *Proc. 2025 American Control Conference*, 2025.

[24] W. Messner and D. Tilbury, *Control Tutorials for MATLAB and Simulink: A Web-Based Approach*. Prentice Hall, 1999.

[25] D. Tilbury and B. Messner, "Control tutorials for MATLAB and Simulink," 2025, accessed: 16th July, 2025, [Online.] Available: https://ctms.engin.umich.edu/CTMS/index.php?aux=Home.

[26] M. Blanke, M. Kinnaert, J. Lunze, and M. Staroswiecki, *Diagnosis and Fault-Tolerant Control*, 3rd ed. Springer, 2016.

[27] J. Milošević, H. Sandberg, and K. H. Johansson, "Estimating the impact of cyber-attack strategies for stochastic networked control systems," *IEEE Trans. Control Netw. Syst.*, vol. 7, no. 2, pp. 747–757, 2019.

[28] H. Sasahara, "Adversarial attacks to direct data-driven control for destabilization," in *Proc. IEEE 62nd Conference on Decision and Control (CDC)*, 2023.

[29] T. Kaminaga and H. Sasahara, "Adversarial attack using projected gradient method to direct data-driven control," in *Proc. IEEE Conference on Control Technology and Applications (CCTA)*, 2024, pp. 236–241.

[30] T. Chen and B. A. Francis, *Optimal Sampled-data Control Systems*. Springer, 2012.

[31] H. J. Van Waarde, J. Eising, H. L. Trentelman, and M. K. Camlibel, "Data informativity: A new perspective on data-driven analysis and control," *IEEE Trans. Autom. Control*, vol. 65, no. 11, pp. 4753–4768, 2020.

[32] J. C. Willems, P. Rapisarda, I. Markovsky, and B. L. De Moor, "A note on persistency of excitation," *Systems & Control Letters*, vol. 54, no. 4, pp. 325–329, 2005.

[33] B. Biggio, B. Nelson, and P. Laskov, "Poisoning attacks against support vector machines," in *Proc. 29th International Coference on Machine Learning*, 2012, p. 1467–1474.

[34] A. Kosugi, K. Teranishi, and K. Kogiso, "Experimental validation of the attack-detection capability of encrypted control systems using man-in-the-middle attacks," *IEEE Access*, vol. 12, pp. 10 535–10 547, 2024.

[35] L. Garcia, F. Brasser, M. H. Cintuglu, A.-R. Sadeghi, O. A. Mohammed, and S. A. Zonouz, "Hey, my malware knows physics! Attacking PLCs with physical model aware rootkit." in *Proc. Network and Distributed System Security (NDSS) Symposium*, 2017.

[36] D. P. Bertsekas, *Nonlinear Programming*, 3rd ed. Athena Scientific, 2016.

[37] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization." *Journal of Machine Learning Research*, vol. 12, no. 7, 2011.

[38] R. L. Burden, J. D. Faires, and A. M. Burden, *Numerical Analysis*, 10th ed. Cengage learning, 2015.

[39] M. Blondel, Q. Berthet, M. Cuturi, R. Frostig, S. Hoyer, F. Llinares-López, F. Pedregosa, and J.-P. Vert, "Efficient and modular implicit differentiation," *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, vol. 35, pp. 5230–5242, 2022.

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

[41] M. Xu, T. L. Molloy, and S. Gould, "Revisiting implicit differentiation for learning problems in optimal control," *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, vol. 36, 2024.

[42] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.

[43] D. Duvenaud, J. Z. Kolter, and M. Johnson, "Deep implicit layers tutorial-neural ODEs, deep equilibirum models, and beyond," *Neural Information Processing Systems Tutorial*, 2020.

[44] J. R. Magnus and H. Neudecker, *Matrix Differential Calculus with Applications in Statistics and Econometrics*. John Wiley & Sons, 2019.

[45] X. Yuan, P. He, Q. Zhu, and X. Li, "Adversarial examples: Attacks and defenses for deep learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 9, pp. 2805–2824, 2019.

[46] MathWorks, "MATLAB version 24.2.0 (R2024b)," Software, Natick, MA, USA, 2024, https://www.mathworks.com.

[47] D. Bernstein, *Scalar, Vector, and Matrix Mathematics: Theory, Facts, and Formulas-revised and Expanded Edition*. Princeton University Press, 2018.

**Hampei Sasahara** (M'19) is Assistant Professor with the Department of Systems and Control Engineering, Institute of Science Tokyo, Tokyo, Japan. He received the Ph.D. degree in engineering from Tokyo Institute of Technology in 2019. From 2019 to 2021, he was a Postdoctoral Scholar with KTH Royal Institute of Technology, Stockholm, Sweden. From 2022 to 2024, he was an Assistant Professor with Tokyo Institute of Technology, Tokyo, Japan. His main interests include secure control system design and control of large-scale systems.