

# From Passive Monitoring to Active Defence: Resilient Control of Manipulators Under Cyberattacks

Gabriele Gualandi and Alessandro V. Papadopoulos

**Abstract**—Cyber-physical robotic systems are vulnerable to *false data injection attacks* (FDIAs), in which an adversary corrupts sensor signals while evading residual-based *passive* anomaly detectors such as the  $\chi^2$  test. Such *stealthy* attacks can induce substantial end-effector deviations without triggering alarms. This paper studies the resilience of redundant manipulators to stealthy FDIAs and advances the architecture from passive monitoring to active defence. We formulate a closed-loop model comprising a feedback-linearized manipulator, a steady-state Kalman filter, and a  $\chi^2$ -based anomaly detector. Building on this passive monitoring layer, we propose an active control-level defence that attenuates the control input through a monotone function of an anomaly score generated by a novel actuation-projected, measurement-free state predictor. The proposed design provides probabilistic guarantees on nominal actuation loss and preserves closed-loop stability. From the attacker perspective, we derive a convex QCQP for computing one-step optimal stealthy attacks. Simulations on a 6-DOF planar manipulator show that the proposed defence significantly reduces attack-induced end-effector deviation while preserving nominal task performance in the absence of attacks.

## I. INTRODUCTION

Robotic manipulators are increasingly deployed in open and networked environments, ranging from industrial assembly to collaborative human-robot interaction. Their tight integration of computation, communication, and actuation, however, makes them vulnerable to cyberattacks that directly compromise safety and reliability. Among cyber-physical threats, attacks on *data integrity* are particularly critical: by corrupting sensor information, an adversary can mislead the controller and alter the robot’s behavior without any physical contact [1], [2]. Such threats are especially concerning when they are *stealthy*, i.e., engineered to remain below the threshold of an anomaly detection system (ADS), thereby avoiding alarms while still driving the end-effector away from its intended task [3], [4].

A well-studied class of integrity attacks are *False Data Injection Attacks* (FDIAs), in which adversaries manipulate sensor signals to achieve malicious objectives. In networked control and power systems, stealthy FDIAs have been extensively analyzed: from characterizations of undetectability [5], [6], [7] to optimal attack synthesis. In robotics, however, prior work has primarily focused on passive anomaly detection or on “perfectly undetectable” attacks that completely bypass detection [7]. These perspectives miss a key robotics-specific vulnerability: the widespread use of *feedback lin-*

*earization* reduces manipulator dynamics to double integrators, which in turn induces an *integrator vulnerability*. As a consequence, persistent sensor corruption can silently accumulate in the closed-loop system, driving large task-space errors even under residual-based  $\chi^2$  detection [8].

**This paper addresses this gap by introducing an active defence strategy that transforms anomaly detection from passive monitoring into a resilience mechanism.** The proposed method leverages an *actuation-projected state predictor*, a model-driven state estimate that ignores sensing, to compute an anomaly score immune to direct sensor corruption. Based on this score, we introduce *anomaly-aware command scaling*, which attenuates commanded accelerations as anomalies grow, thereby reducing the adversary’s ability to steer the manipulator while preserving nominal performance.

*Contributions.*: The contributions of this work are the following:

- 1) We formalize stealthy false data injection attacks (FDIAs) against the sensors of feedback-linearized manipulators, exposing their integrator vulnerability and showing that the attacker’s one-step optimal strategy reduces to a convex QCQP.
- 2) We propose *anomaly-aware command scaling*, which attenuates control inputs based on a measurement-free actuation-projected state predictor.
- 3) We provide probabilistic guarantees on bounded attenuation in nominal operation and prove closed-loop stability under the proposed defence.
- 4) Simulations on a 6-DOF redundant manipulator demonstrate that the defence significantly limits attacker-realizable end-effector deviations while preserving task performance in the absence of attacks.

## A. Related Work

**Anomaly detection in CPS.** Residual-based detectors, particularly  $\chi^2$  tests on Kalman innovations, are a classical tool for monitoring CPS integrity [9]. These methods provide statistical guarantees on false-alarm rates and are widely used in industrial practice. Extensions include adaptive thresholds [10] and sequential schemes such as CUSUM tests [11]. However, such ADS are inherently *passive*: they detect anomalies but do not alter the control policy, leaving the system structure unchanged.

**False Data Injection Attacks.** Theoretical studies of FDIAs in networked systems have characterized undetectability conditions [5], [6], affine attack structures [7], and optimal attack strategies [3]. These works typically assume

This work was supported by the Swedish Research Council (VR) with the PSI project No. #2020-05094, and by the Knowledge Foundation (KKS).

G. Gualandi and A.V. Papadopoulos are with the Department of Computer Science and Engineering, Mälardalen University, Västerås, Sweden [gabriele.gualandi@mdu.se](mailto:gabriele.gualandi@mdu.se)

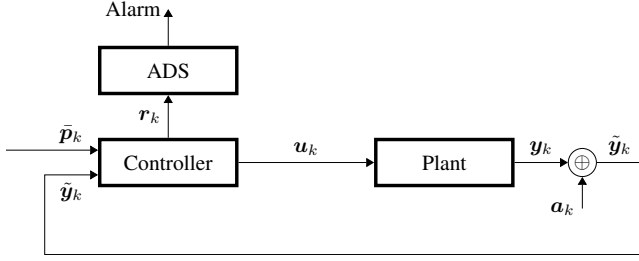


Fig. 1. Closed-loop system under FDIA. Sensor signals are corrupted by injection  $\mathbf{a}_k$ , yielding attacked output  $\tilde{\mathbf{y}}_k$ . The Kalman filter generates the innovation  $\mathbf{r}_k$ , which is monitored by the ADS.

either “perfectly undetectable” attacks (no residual information leaks) or focus on power networks and generic LTI plants. Robotics-specific studies remain limited: most works concentrate on detector design rather than on modifying the controller to actively limit attack effectiveness [4].

**Integrator Vulnerability.** Control systems with an LTI plant having integral action are subject to integrator vulnerability [8], where the residual generated by any linear observer follows the same distribution during normal operation and under attack in the steady-state regime, making sensor bias injecting attack detectable only during transients.

**Our contribution.** We bridge these lines of research by combining a residual-based ADS with a novel active defence operating at the control level. First, we show that feedback linearization induces an *integrator vulnerability*: a PD-based FDIA can precisely steer the end-effector despite a residual-based  $\chi^2$  detector. We then show that our gain-scaling defence, driven by the anomaly score, reduces the attack’s effectiveness while providing probabilistic guarantees in attack-free operation, thereby improving the cyber-physical security of robotic manipulators.

## II. SYSTEM MODEL

We consider a robotic manipulator operating in closed loop with a state estimator, a task-space controller, and an anomaly detection system (ADS). The architecture, illustrated in Fig. 1, captures the interaction between defender and adversary: (i) the plant dynamics, (ii) a Kalman filter for state estimation, (iii) task-space control, (iv) a residual-based ADS, and (v) an additive adversarial attack on sensor measurements.

### A. Closed-Loop Architecture

At each discrete time step  $k \in \mathbb{Z}_{\geq 0}$ , the plant output  $\mathbf{y}_k \in \mathbb{R}^p$  is corrupted by an injected signal  $\mathbf{a}_k$  to yield the attacked measurement

$$\tilde{\mathbf{y}}_k = \mathbf{y}_k + \mathbf{a}_k, \quad (1)$$

### B. Plant and State Estimator

The plant is modeled as a discrete-time LTI system:

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{w}_k, \quad (2a)$$

$$\mathbf{y}_k = \mathbf{C}\mathbf{x}_k + \mathbf{v}_k, \quad (2b)$$

with state  $\mathbf{x}_k \in \mathbb{R}^n$ , process noise  $\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$ , and measurement noise  $\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$ . We assume  $(\mathbf{A}, \mathbf{C})$  is detectable, ensuring the existence of a unique stabilizing solution  $\mathbf{P}$  to the DARE

$$\mathbf{P} = \mathbf{A}\mathbf{P}\mathbf{A}^\top + \mathbf{Q} - \mathbf{A}\mathbf{P}\mathbf{C}^\top (\mathbf{C}\mathbf{P}\mathbf{C}^\top + \mathbf{R})^{-1} \mathbf{C}\mathbf{P}\mathbf{A}^\top.$$

State estimates are obtained from a steady-state Kalman filter in single-step innovation form:

$$\hat{\mathbf{x}}_{k+1} = \mathbf{A}\hat{\mathbf{x}}_k + \mathbf{B}\mathbf{u}_k + \mathbf{L}\mathbf{r}_k, \quad (3a)$$

$$\mathbf{r}_k = \tilde{\mathbf{y}}_k - \mathbf{C}\hat{\mathbf{x}}_k, \quad (3b)$$

where  $\mathbf{r}_k \in \mathbb{R}^p$  is the innovation. Its steady-state covariance is  $\Sigma = \mathbf{C}\mathbf{P}\mathbf{C}^\top + \mathbf{R}$ , and the corresponding steady-state gain is

$$\mathbf{L} = \mathbf{A}\mathbf{P}\mathbf{C}^\top \Sigma^{-1}. \quad (4)$$

The residual  $\mathbf{r}_k$  and its covariance  $\Sigma$  will later be used for residual-based anomaly detection.

### C. Manipulator Model and Task-Space Control

We consider a  $p$ -DOF robotic manipulator with sensor readings providing joint values  $\mathbf{q}$ , and having state  $\mathbf{x} = [\mathbf{q}^\top \dot{\mathbf{q}}^\top]^\top \in \mathbb{R}^n$ , where  $n = 2 \cdot p$  and  $\mathbf{q} \in \mathbb{R}^p$ .

The continuous-time joint-space dynamics are

$$\mathbf{M}(\mathbf{q})\ddot{\mathbf{q}} + \boldsymbol{\nu}(\mathbf{q}, \dot{\mathbf{q}}) = \boldsymbol{\tau}, \quad (5)$$

where  $\mathbf{M}(\mathbf{q}) \succ 0$  is the inertia matrix and  $\boldsymbol{\nu}$  collects Coriolis, centrifugal, and gravity terms. With inverse-dynamics compensation, i.e.,  $\boldsymbol{\tau} = \mathbf{M}(\mathbf{q})\mathbf{u}^{\text{nom}} + \boldsymbol{\nu}(\mathbf{q}, \dot{\mathbf{q}})$ , the joint dynamics reduce to decoupled double integrators,

$$\ddot{\mathbf{q}} = \mathbf{u}^{\text{nom}}, \quad (6)$$

whose discrete-time form is represented by Eq. (2) with  $\mathbf{x}_k = [\mathbf{q}_k^\top \dot{\mathbf{q}}_k^\top]^\top$  and input  $\mathbf{u}_k = \mathbf{u}_k^{\text{nom}}$ .

The controller operates in task space using twist control [12], [13], tracking position references  $\{\bar{\mathbf{p}}_k, \dot{\bar{\mathbf{p}}}_k, \ddot{\bar{\mathbf{p}}}_k\}$  and orientation references  $\{\bar{\mathbf{R}}_k, \bar{\boldsymbol{\omega}}_k, \dot{\bar{\boldsymbol{\omega}}}_k\}$ , where  $\bar{\mathbf{R}}_k \in SO(3)$  and  $\bar{\boldsymbol{\omega}}_k \in \mathbb{R}^3$ .

Position and orientation errors are defined as

$$\mathbf{e}_{p,k} = \bar{\mathbf{p}}_k - \hat{\mathbf{p}}_k, \quad (7)$$

$$\mathbf{e}_{o,k} = \sin\left(\frac{\hat{\theta}_k}{2}\right) \hat{\mathbf{r}}_k, \quad (8)$$

where  $(\hat{\theta}_k, \hat{\mathbf{r}}_k)$  is the angle-axis pair of the rotation error  $\bar{\mathbf{R}}_k \hat{\mathbf{R}}_k^\top$  with  $\hat{\theta}_k \in [0, \pi]$ . The estimates  $\hat{\mathbf{p}}_k$  and  $\hat{\mathbf{R}}_k$  are obtained via forward kinematics from  $\hat{\mathbf{q}}_k$  extracted from the state estimate  $\hat{\mathbf{x}}_k$ .

A PD+feedforward law generates desired task accelerations:

$$\mathbf{u}_k^c = \begin{bmatrix} \ddot{\bar{\mathbf{p}}}_k + \mathbf{K}_{pp} \mathbf{e}_{p,k} + \mathbf{K}_{dp} \dot{\mathbf{e}}_{p,k} \\ \dot{\bar{\boldsymbol{\omega}}}_k + \mathbf{K}_{po} \mathbf{e}_{o,k} + \mathbf{K}_{do} \dot{\mathbf{e}}_{o,k} \end{bmatrix}, \quad (9)$$

where  $\dot{\mathbf{e}}_{p,k} := \dot{\bar{\mathbf{p}}}_k - \dot{\hat{\mathbf{p}}}_k$  and  $\dot{\mathbf{e}}_{o,k} := \dot{\bar{\boldsymbol{\omega}}}_k - \dot{\hat{\boldsymbol{\omega}}}_k$  are velocity errors. The gains are synthesized via discrete-time LQR on the discretized double-integrator model (6).

Joint accelerations are then computed via the pseudoinverse of the geometric Jacobian  $\mathbf{J}(\mathbf{q}) \in \mathbb{R}^{6 \times p}$  [13]:

$$\mathbf{u}_k^{\text{nom}} = \mathbf{J}^\dagger(\mathbf{q})(\mathbf{u}_k^c - \dot{\mathbf{J}}(\mathbf{q}, \dot{\mathbf{q}})\dot{\mathbf{q}}). \quad (10)$$

#### D. Passive Anomaly Detection System

The ADS monitors the innovation sequence  $\mathbf{r}_k$ . With steady-state covariance  $\Sigma \succ 0$  (hence invertible), the Mahalanobis distance

$$z_k = \mathbf{r}_k^\top \Sigma^{-1} \mathbf{r}_k \quad (11)$$

is  $\chi^2(p)$ -distributed under  $\mathcal{H}_0$  (no attack), where  $p = \dim(\mathbf{y}_k)$ . To reduce sensitivity to single-sample fluctuations, we employ a windowed statistic of length  $W \in \mathbb{N}$ :

$$w_k = \sum_{i=k-W+1}^k z_i. \quad (12)$$

During execution, an alarm is triggered if  $w_k > \tau$ .

*Lemma 1 (Innovation whiteness):* Under  $\mathcal{H}_0$  (no attack), the innovation process satisfies

$$\mathbf{r}_k \sim \mathcal{N}(\mathbf{0}, \Sigma), \quad \mathbf{r}_k \text{ independent across } k.$$

*Proof:* This is a standard property of the Kalman filter with correct noise statistics: the innovation sequence is white, zero-mean Gaussian with covariance  $\Sigma$ , see [9]. ■

*Lemma 2 (Chi-squared distribution):* Under  $\mathcal{H}_0$ , the statistics satisfy

$$z_k \sim \chi^2(p), \quad w_k \sim \chi^2(pW),$$

where  $p$  is the output dimension and  $W$  the window length.

*Proof:* From Lemma 1,  $\Sigma^{-1/2} \mathbf{r}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . Thus  $z_k = \|\Sigma^{-1/2} \mathbf{r}_k\|_2^2 \sim \chi^2(p)$  [14]. Independence across time implies  $w_k$  is the sum of  $W$  i.i.d.  $\chi^2(p)$  variables, hence  $\chi^2(pW)$ , see [10]. ■

*Corollary 1 (Threshold calibration):* Given a desired per-step false-alarm probability  $\alpha \in (0, 1)$ , set

$$\tau = F_{\chi^2(pW)}^{-1}(1 - \alpha) = 2P^{-1} \left( 1 - \alpha, \frac{pW}{2} \right), \quad (13)$$

which ensures  $\Pr(w_k > \tau \mid \mathcal{H}_0) = \alpha$ , where  $P^{-1}$  denotes the inverse of the regularized lower incomplete gamma function [10].

#### E. Attacker Model

We now state the assumptions under which both the defence and attack strategies are developed.

*Assumption 1 (Adversary knowledge):* The adversary is omniscient: it has full knowledge of the plant, controller, state estimator, and defence system(s).

*Assumption 2 (Adversary capabilities):* The only attack surface is additive injection into sensor measurements. The adversary cannot tamper with actuators, control logic, ADS parameters, timing, or packet ordering.

*Assumption 3 (Adversary goal):* The adversary's goal is to manipulate the position of the end-effector in task space, as formalized later by an attack objective on the task-space trajectory, subject to the detection mechanism.

*Assumption 4 (Converged estimator):* Without loss of generality, the attack starts at  $k = 0$  and lasts at most  $T$  samples. At  $k = 0$ , the Kalman filter operates at steady state (covariances and gains constant), i.e., the estimator has converged.

### III. PROPOSED DEFENCE METHOD

We now present a defence strategy based on gain scaling to reduce the impact of sensing FDIA, with probabilistic guarantees on the performance loss under attack-free operation.

#### A. Attack Estimation

We define the *actuation-projected state*  $\tilde{\mathbf{x}}_k$  as the measurement-free state prediction driven solely by commanded inputs. To mitigate drift from model mismatch, this predictor is periodically re-synchronized with the Kalman estimate as ground truth. Without loss of generality, assume the most recent re-synchronization occurs at  $k = 0$ , so

$$\tilde{\mathbf{x}}_0 = \hat{\mathbf{x}}_0. \quad (14)$$

For subsequent steps, the actuation-projected state follows the noise-free, open-loop predictor

$$\tilde{\mathbf{x}}_{k+1} = \mathbf{A} \tilde{\mathbf{x}}_k + \mathbf{B} \mathbf{u}_k, \quad k \geq 0. \quad (15)$$

We define the *actuation-projected residual* (distinct from the innovation  $\mathbf{r}_k$ ) as

$$\tilde{\mathbf{r}}_k := \hat{\mathbf{x}}_k - \tilde{\mathbf{x}}_k \in \mathbb{R}^n, \quad n = 2p. \quad (16)$$

*Theorem 1 (Actuation-projected residual covariance):*

The covariance of the residual in Eq. (16) after  $k$  steps from initialization,  $\Sigma_{\tilde{\mathbf{r}},k} \in \mathbb{R}^{n \times n}$ , is the  $(2, 2)$  block of the matrix  $\mathbf{P}_{z,k} \in \mathbb{R}^{2n \times 2n}$  obtained by iterating

$$\mathbf{P}_{z,j+1} = \mathbf{F} \mathbf{P}_{z,j} \mathbf{F}^\top + \mathbf{\Pi}, \quad j = 0, \dots, k-1, \quad (17)$$

from the initial condition  $\mathbf{P}_{z,0} = \text{diag}(\mathbf{P}, \mathbf{0})$ , where

$$\mathbf{F} = \begin{bmatrix} \mathbf{A} - \mathbf{LC} & \mathbf{0} \\ \mathbf{LC} & \mathbf{A} \end{bmatrix} \in \mathbb{R}^{2n \times 2n}, \quad (18)$$

$$\mathbf{G} = \begin{bmatrix} \mathbf{I} & -\mathbf{L} \\ \mathbf{0} & \mathbf{L} \end{bmatrix} \in \mathbb{R}^{2n \times (n+p)}, \quad (19)$$

$$\mathbf{\Pi} = \mathbf{G} \text{diag}(\mathbf{Q}, \mathbf{R}) \mathbf{G}^\top \in \mathbb{R}^{2n \times 2n}. \quad (20)$$

*Proof:* From the system equations, the one-step evolution of the estimation error  $\mathbf{e}_k := \mathbf{x}_k - \hat{\mathbf{x}}_k$  and of the actuation-projected residual  $\tilde{\mathbf{r}}_k := \hat{\mathbf{x}}_k - \tilde{\mathbf{x}}_k$  is:

$$\mathbf{e}_{k+1} = (\mathbf{A} - \mathbf{LC}) \mathbf{e}_k + \mathbf{w}_k - \mathbf{L} \mathbf{v}_k, \quad (21)$$

$$\tilde{\mathbf{r}}_{k+1} = \mathbf{A} \tilde{\mathbf{r}}_k + \mathbf{LC} \mathbf{e}_k + \mathbf{L} \mathbf{v}_k. \quad (22)$$

Stacking  $\mathbf{z}_k := [\mathbf{e}_k^\top, \tilde{\mathbf{r}}_k^\top]^\top$  yields

$$\mathbf{z}_{k+1} = \mathbf{F} \mathbf{z}_k + \mathbf{G} \boldsymbol{\eta}_k, \quad (23)$$

with  $\mathbf{F}, \mathbf{G}$  as in Eqs. (18) and (19) and  $\boldsymbol{\eta}_k \triangleq [\mathbf{w}_k^\top, \mathbf{v}_k^\top]^\top$ . The covariance propagates as

$$\begin{aligned} \mathbf{P}_{z,k+1} &= \mathbf{F} \mathbb{E}[\mathbf{z}_k \mathbf{z}_k^\top] \mathbf{F}^\top + \mathbf{F} \mathbb{E}[\mathbf{z}_k \boldsymbol{\eta}_k^\top] \mathbf{G}^\top \\ &\quad + \mathbf{G} \mathbb{E}[\boldsymbol{\eta}_k \mathbf{z}_k^\top] \mathbf{F}^\top + \mathbf{G} \mathbb{E}[\boldsymbol{\eta}_k \boldsymbol{\eta}_k^\top] \mathbf{G}^\top. \end{aligned}$$

Since  $\boldsymbol{\eta}_k$  is white, zero-mean, and independent of  $\mathbf{z}_k$ ,  $\mathbb{E}[\mathbf{z}_k \boldsymbol{\eta}_k^\top] = \mathbf{0}$  and  $\mathbb{E}[\boldsymbol{\eta}_k \boldsymbol{\eta}_k^\top] = \text{diag}(\mathbf{Q}, \mathbf{R})$ , giving

$$\mathbf{P}_{z,k+1} = \mathbf{F} \mathbf{P}_{z,k} \mathbf{F}^\top + \mathbf{\Pi},$$

with  $\mathbf{\Pi}$  as in Eq. (20).

At synchronization  $k = 0$ ,  $\text{Cov}(\mathbf{e}_0) = \mathbf{P}$  and  $\tilde{\mathbf{r}}_0 = \mathbf{0}$ , hence

$$\mathbf{P}_{z,0} = \begin{bmatrix} \mathbf{P} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} = \text{diag}(\mathbf{P}, \mathbf{0}).$$

Iterating (17)  $k$  times yields  $\mathbf{P}_{z,k}$ ; the desired covariance is the (2, 2) block,  $\Sigma_{\tilde{\mathbf{r}},k}$ . ■

*Theorem 2 (Confidence for the Anomaly Measure):*

Consider the anomaly measure

$$\tilde{z}_k = \tilde{\mathbf{r}}_k^\top \Sigma_{\tilde{\mathbf{r}},k}^{-1} \tilde{\mathbf{r}}_k. \quad (24)$$

Choosing  $z_x = F_{\chi^2(n)}^{-1}(\psi)$  ensures  $\mathbb{P}(\tilde{z}_k \leq z_x \mid \mathcal{H}_0) = \psi$  for any given probability  $\psi$ .

*Proof:* Under  $\mathcal{H}_0$ , the residual  $\tilde{\mathbf{r}}_k$  is a zero-mean Gaussian vector,  $\tilde{\mathbf{r}}_k \sim \mathcal{N}(\mathbf{0}, \Sigma_{\tilde{\mathbf{r}},k})$ . The time-varying normalization with the covariance  $\Sigma_{\tilde{\mathbf{r}},k}$  at each step, computed as in Theorem 1, ensures that the resulting statistic has a stationary distribution,  $\tilde{z}_k \sim \chi^2(n)$ , for all  $k$ . The theorem's result follows from the definition of the inverse CDF,  $F_{\chi^2(n)}^{-1}$ . ■

### B. Command scaling active defence

As an active defence mechanism, we propose to reduce the magnitude of command  $\mathbf{u}_k^{\text{nom}}$  of Eq. (10) through a gain scaling function  $f(\tilde{z})$ , with  $\tilde{z}$  as in Eq. (24).

Let  $z_x > 0$  be a design abscissa and  $\beta \in (0, 1)$  the desired gain factor at  $\tilde{z} = z_x$ . Introduce a shape exponent  $\gamma > 0$  and a smooth, strictly decreasing function  $f : [0, \infty) \rightarrow (0, 1]$ ,

$$f(\tilde{z}) = \exp \left[ - \left( \frac{\tilde{z}}{z_{\text{scale}}} \right)^\gamma \right], \quad (25)$$

$$z_{\text{scale}} = \frac{z_x}{(-\ln(\beta))^{1/\gamma}}, \quad (26)$$

satisfying  $f(0) = 1$ ,  $f(z_x) = \beta$ , and  $f(\tilde{z}) \rightarrow 0$  as  $\tilde{z} \rightarrow \infty$ .

*Theorem 3:* Consider a robotic manipulator with nominal joint-space control input  $\mathbf{u}_k^{\text{nom}}$  as in (10). Let  $\tilde{z}_k \in [0, \infty)$  denote the anomaly score defined in (24), and let  $f : [0, \infty) \rightarrow (0, 1]$  be a strictly decreasing function satisfying

$$f(0) = 1, \quad \lim_{\tilde{z} \rightarrow \infty} f(\tilde{z}) = 0.$$

Define the scaled control law

$$\mathbf{u}_k = f(\tilde{z}_k) \mathbf{u}_k^{\text{nom}}. \quad (27)$$

Then the following hold:

- 1) **Probabilistic actuation guarantee.** Let  $\tilde{z}_k \sim \chi^2(n)$  under the null hypothesis, and let  $z_x = F_{\chi^2(n)}^{-1}(\psi)$  denote the inverse CDF at confidence level  $\psi \in (0, 1)$ . If design parameters of  $f(\cdot)$  are chosen such that

$$f(z_x) \geq \beta,$$

for some  $\beta \in (0, 1)$ , then with probability at least  $\psi$ , the scaled input magnitude satisfies

$$\|\mathbf{u}_k\| \geq \beta \|\mathbf{u}_k^{\text{nom}}\|. \quad (28)$$

- 2) **Closed-loop stability (frozen gain).** Consider the discretized double-integrator dynamics induced by inverse-dynamics compensation in (6). If the nominal frozen-gain closed loop is Schur stable (i.e., (27) with  $f(\cdot) \equiv 1$ ), then for any constant gain factor  $\bar{f} \in (0, 1]$  the frozen-gain closed loop under Eq. (27) is Schur stable, and hence exponentially stable.

*Proof: Probabilistic actuation guarantee.* Under  $\mathcal{H}_0$ , the anomaly score follows  $\tilde{z}_k \sim \chi^2(n)$ . By definition of the quantile,  $\Pr(\tilde{z}_k \leq z_x) = \psi$ . On this event, since  $f(\cdot)$  is non-increasing, we have  $f(\tilde{z}_k) \geq f(z_x)$ . By assumption,  $f(z_x) \geq \beta$ . Substituting into (27),

$$\|\mathbf{u}_k\| = f(\tilde{z}_k) \|\mathbf{u}_k^{\text{nom}}\| \geq \beta \|\mathbf{u}_k^{\text{nom}}\|.$$

Therefore, with probability at least  $\psi$ , the commanded input magnitude remains at least a fraction  $\beta$  of the nominal value.

**Closed-loop stability (frozen gain).** Under inverse-dynamics compensation, and treating the Jacobian pseudoinverse as locally constant, it suffices to consider one joint with  $\mathbf{u}_{k,j}^{\text{nom}} = -\mathbf{K}\mathbf{x}_{k,j}$ , where  $\mathbf{x}_{k,j} \in \mathbb{R}^2$  and  $\mathbf{K} = [K_p \ K_d]$ . For a frozen  $\bar{f} \in (0, 1]$ , the closed-loop matrix is  $\mathbf{A}_{\text{cl}}(\bar{f}) = \mathbf{A}_j - \bar{f}\mathbf{B}_j\mathbf{K}$ , where  $\mathbf{A}_j, \mathbf{B}_j$  are the joint double-integrator matrices. Applying the second-order Jury criterion to the characteristic polynomial of  $\mathbf{A}_{\text{cl}}(\bar{f})$ , Schur stability is equivalent to  $\bar{f}T_s^2K_p > 0$ ,  $2 - \bar{f}T_sK_d > 0$ , and  $K_d > \frac{T_s}{2}K_p$ . These conditions hold at  $\bar{f} = 1$  by assumption, hence  $K_p > 0$ ,  $2 > T_sK_d$ , and  $K_d > \frac{T_s}{2}K_p$ . Therefore, for every  $\bar{f} \in (0, 1]$ , the first condition holds since  $\bar{f} > 0$  and  $K_p > 0$ , the second becomes less restrictive as  $\bar{f}$  decreases, and the third is independent of  $\bar{f}$ . Thus  $\mathbf{A}_{\text{cl}}(\bar{f})$  is Schur for all  $\bar{f} \in (0, 1]$ , and the frozen-gain closed loop is exponentially stable. ■

In practice,  $f(\tilde{z}_k)$  is time-varying and state-dependent; however, item 1 of the theorem confines  $f$  to the compact set  $[\beta, 1]$  with probability  $\psi$ , over which frozen stability holds uniformly.

*Corollary 2:* Fix a desired confidence level  $\psi \in (0, 1)$  and a desired gain floor  $\beta \in (0, 1)$ . Fix  $z_x = F_{\chi^2(n)}^{-1}(\psi)$  in (25) so that  $f(z_x) = \beta$ . Then, for the control law (27) with  $\tilde{z}_k$  as in Eq. (24), item 1 of Theorem 3 guarantees that the magnitude of the final command satisfies (28) with probability  $\psi$ , i.e., the final input remains at least  $\beta \cdot 100\%$  of the nominal magnitude with probability  $\psi$ .

## IV. OPTIMAL STEALTH ATTACK

We now formalize the adversary's strategy under Assumptions of Section II-E. The attacker's goal is to induce stealthy malicious end-effector accelerations.

To express the predicted trajectories of different vectors (e.g., acceleration, velocity) from a closed-loop simulation, we introduce the following compact notation.

Let  $\mathcal{Z}_{k,j}(\mathbf{a})$  denote the  $j$ -step-ahead prediction at time  $k+j$ , computed at time  $k$ , of the closed-loop system under the attack sequence  $[\mathbf{a}, \mathbf{0}, \dots, \mathbf{0}]$  (of length  $j+1$ ).

For a generic vector  $\mathbf{v}$ , we define

$$\mathbf{v}_{k+j}^{\text{SIM}} := \pi_{\mathbf{v}}(\mathcal{Z}_{k,j}(\mathbf{a})), \quad (29)$$

where  $\pi_v(\cdot)$  extracts  $v$  from the noise-free simulated trajectory at prediction step  $j$  (i.e., the terminal value at time  $k+j$ ). The operator  $\mathcal{Z}_{k,j}(\mathbf{a})$  models Assumption 1 of Section II-E, as it propagates Eqs. (2), (3a) and (3b) under perfect-model assumptions (exact initial state and noise-free dynamics) and it computes the control input via Eq. (27) including the proposed active defence mechanism.

#### A. Attack Model and Delay Structure

Due to estimator and control delays, an injected signal at time  $k$ ,  $\mathbf{a}_k$ , influences the end-effector acceleration two steps later (under the closed-loop dynamics in Eqs. (2), (3a) and (3b)). Specifically,

$$\ddot{\mathbf{p}}_{k+2}^{\text{SIM}} = \pi_{\ddot{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_k)), \quad (30)$$

where  $\mathcal{Z}_{k,2}$  denotes a two-step-ahead simulation of the closed-loop system under attack sequence  $[\mathbf{a}_k, \mathbf{0}, \mathbf{0}]$ .

The attacker's high-level objective at each time step is to make the predicted end-effector acceleration match a desired target acceleration  $\ddot{\mathbf{p}}_{k+2}^{\text{A}}$ . This is formulated as:

$$\min_{\mathbf{a}_k} \frac{1}{2} \left\| \ddot{\mathbf{p}}_{k+2}^{\text{A}} - \pi_{\ddot{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_k)) \right\|^2. \quad (31)$$

#### B. Incremental Attack Formulation

Due to feedback linearization, the manipulator's joint-space dynamics reduce to double integrators, which entails an *integrator vulnerability*: a persistent sensor bias can induce drift in the regulated variables. This vulnerability has been analyzed for sensor *bias injection attacks* (BIAs)—i.e., constant sensor injections—in linear systems [8]. Although the joint-task mapping is nonlinear and we consider more general *false data injection attacks* (FDIAs) with arbitrary sensor injections, we argue that effective FDIAs still exploit the integrator vulnerability *locally*. For this reason, we model the attack *incrementally*:

$$\mathbf{a}_k = \mathbf{a}_{k-1} + \Delta_k, \quad k \geq 1, \quad (32)$$

where  $\Delta_k$  is an increment and  $\mathbf{a}_k$  is initialized as  $\mathbf{a}_0 = \Delta_0$  at the attack onset. This parameterizes an FDIA as deviations from a baseline BIA (via the increments  $\Delta_k$ ), and is convenient for gradient-based synthesis.

A first-order expansion of  $\pi_{\ddot{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}))$  around  $\mathbf{a}_{k-1}$ , with  $\ddot{\mathbf{p}}_{k+2}^{\text{SIM}}$  as in Eq. (30), gives:

$$\ddot{\mathbf{p}}_{k+2}^{\text{SIM}} \approx \pi_{\ddot{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_{k-1})) + \mathbf{Z}_k \Delta_k + \mathcal{O}(\|\Delta_k\|^2), \quad (33)$$

where the Jacobian

$$\mathbf{Z}_k = \left. \frac{\partial}{\partial \mathbf{a}} \pi_{\ddot{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a})) \right|_{\mathbf{a}=\mathbf{a}_{k-1}} \quad (34)$$

is computed numerically via a central-difference scheme with step-size tuning [15];  $\mathbf{Z}_k$  provides the local sensitivity needed to differentiate the objective Eq. (31).

By substituting the linear attack model (33) into (31) and introducing a regularization term, a quadratic objective

function is obtained. Specifically, the attack increment  $\Delta_k$  minimizes the quadratic cost

$$\frac{1}{2} \left\| \mathbf{Z}_k \Delta_k - \ddot{\mathbf{p}}_{k+2}^{\text{A}} + \pi_{\ddot{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_{k-1})) \right\|^2 + \frac{\zeta}{2} \|\Delta_k\|^2, \quad (35)$$

where  $\zeta > 0$  penalizes large increments.

#### C. Attack Objective

We restrict the attack to the translational DOFs, as these are most relevant to adversarial manipulation. The adversary defines  $\ddot{\mathbf{p}}_{k+2}^{\text{A}}$  of Eq. (31) via a one-step-ahead PD law plus feedforward based on a plan  $\{\bar{\mathbf{p}}_k^{\text{A}}, \dot{\bar{\mathbf{p}}}_k^{\text{A}}, \ddot{\bar{\mathbf{p}}}_k^{\text{A}}\}$ . Specifically, the PD errors are computed as the difference between the desired and predicted quantities, under the incremental attack approach ( $\Delta_k=0$  or equivalently,  $\mathbf{a}_k = \mathbf{a}_{k-1}$ ):

$$\begin{aligned} \ddot{\mathbf{p}}_{k+2}^{\text{A}} &= \mathbf{K}_p^{\text{A}} (\bar{\mathbf{p}}_{k+1}^{\text{A}} - \mathbf{p}_{k+1}^{\text{SIM}}) \\ &\quad + \mathbf{K}_d^{\text{A}} (\dot{\bar{\mathbf{p}}}_{k+1}^{\text{A}} - \dot{\mathbf{p}}_{k+1}^{\text{SIM}}) + \ddot{\bar{\mathbf{p}}}_k^{\text{A}}, \end{aligned} \quad (36)$$

where

$$\begin{aligned} \mathbf{p}_{k+1}^{\text{SIM}} &= \pi_{\mathbf{p}}(\mathcal{Z}_{k,1}(\mathbf{a}_{k-1})), \\ \dot{\mathbf{p}}_{k+1}^{\text{SIM}} &= \pi_{\dot{\mathbf{p}}}(\mathcal{Z}_{k,1}(\mathbf{a}_{k-1})), \end{aligned}$$

are the one-step-ahead end-effector position and velocity obtained from the closed-loop simulator under the attack sequence  $[\mathbf{a}_{k-1}, \mathbf{0}]$ . The resulting  $\ddot{\mathbf{p}}_{k+2}^{\text{A}}$  is therefore a way to counteract the predicted drift resulting from  $\Delta_k=0$ .

#### D. Stealth Constraint

The ADS residual is modeled as

$$\mathbf{r}_k = \Delta_k + \mathbf{c}_k, \quad (37)$$

with  $\mathbf{c}_k$  denoting the baseline innovation, defined as

$$\mathbf{c}_k = (\mathbf{y}_k + \mathbf{a}_{k-1} - \hat{\mathbf{y}}_k). \quad (38)$$

The adversary uniformly allocates a per-step anomaly budget

$$\tau' = \frac{\tau}{T}, \quad (39)$$

so that the stealth constraint becomes

$$(\Delta_k + \mathbf{c}_k)^\top \Sigma^{-1} (\Delta_k + \mathbf{c}_k) \leq \tau'. \quad (40)$$

#### E. QCQP Formulation

Considering the objective function Eq. (35) and expanding the constraint of Eq. (40), the adversary's problem reduces to the convex QCQP

$$\begin{aligned} \min_{\Delta_k} \quad & \frac{1}{2} \Delta_k^\top \mathbf{H} \Delta_k + \mathbf{g}^\top \Delta_k \\ \text{s.t.} \quad & \Delta_k^\top \mathbf{O} \Delta_k + \mathbf{b}^\top \Delta_k + c \leq 0, \end{aligned} \quad (41)$$

with parameters

$$\begin{aligned} \mathbf{H} &= \mathbf{Z}_k^\top \mathbf{Z}_k + \zeta \mathbf{I}, \\ \mathbf{g} &= -\mathbf{Z}_k^\top \left( \ddot{\mathbf{p}}_{k+2}^{\text{A}} - \pi_{\ddot{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_{k-1})) \right), \\ \mathbf{O} &= \Sigma^{-1}, \\ \mathbf{b} &= 2\Sigma^{-1} \mathbf{c}_k, \\ c &= \mathbf{c}_k^\top \Sigma^{-1} \mathbf{c}_k - \tau'. \end{aligned}$$

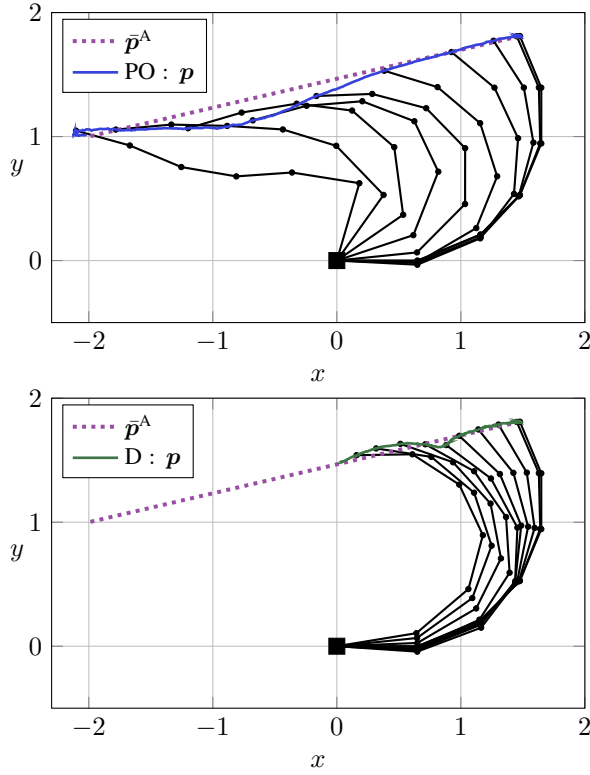


Fig. 2. End-effector trajectories: attacker's reference  $\bar{p}^A$ , defended only by the passive detector (PO) and by the proposed active defence (D). The latter defence limits drift toward the malicious target.

**Theorem 4 (Convexity of adversary's QCQP):** Problem (41) is convex, and thus admits the global minimizer  $\Delta_k^*$ .

*Proof:* The cost Hessian  $\mathbf{H} = \mathbf{Z}_k^\top \mathbf{Z}_k + \zeta \mathbf{I} \succeq \zeta \mathbf{I} > 0$ , so the objective is strictly convex. The constraint has Hessian  $\mathbf{O} = \Sigma^{-1} \succ 0$ , yielding a convex (ellipsoidal) feasible set. A strictly convex objective over a non-empty convex feasible set guarantees existence and uniqueness of the global minimizer  $\Delta_k^*$ . ■

**Remark 1 (Feasibility):** The feasible set is non-empty at every step: the choice  $\Delta_k = -\mathbf{c}_k$  drives the residual to zero and trivially satisfies the constraint for any  $\tau' > 0$ .

#### F. Iterative Attack Construction

The attack evolves incrementally as in Eq. (32)

$$\mathbf{a}_k = \mathbf{a}_{k-1} + \Delta_k^*, \quad (42)$$

initialized with  $\mathbf{a}_0 = \Delta_0^*$ . At each step, the adversary runs an internal simulation to compute  $\mathcal{Z}_k$  and  $\mathbf{Z}_k$ , then solves (41) to determine  $\Delta_k^*$ .

### V. SIMULATION RESULTS

For reproducibility, the Matlab implementation used in this work has been archived and is publicly available [16].

#### A. Experimental Setup

The plant is a 6-DOF planar manipulator with link lengths  $[0.65, 0.55, 0.45, 0.45, 0.45, 0.45]$  m.

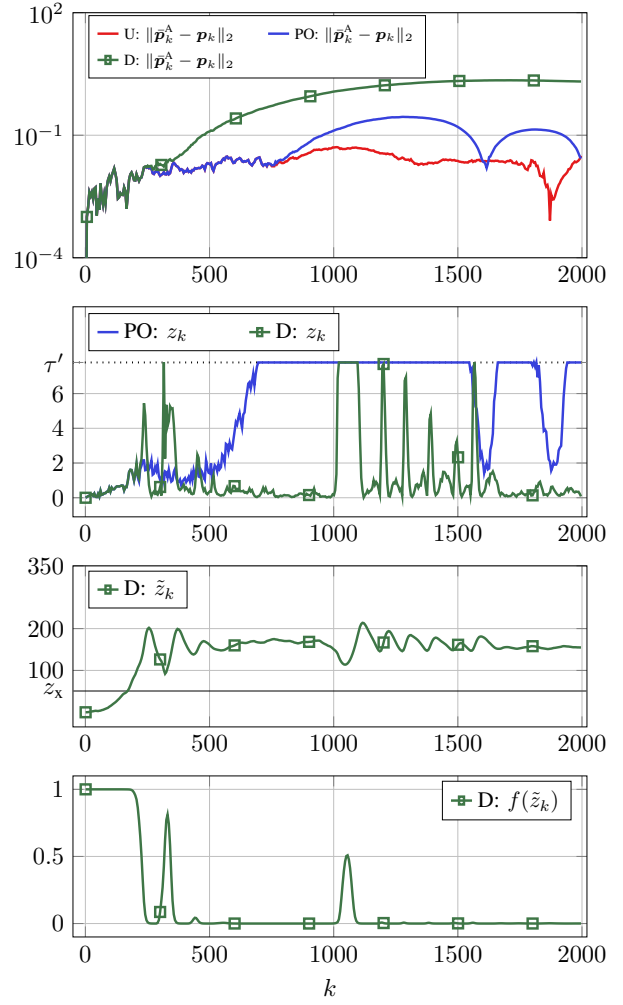


Fig. 3. Attack tracking error  $\|\bar{p}_k^A - \mathbf{p}_k\|_2$ , Mahalanobis distances  $z_k$  and  $\tilde{z}_k$  of Eqs. (11) and (24), and scaling  $f(\tilde{z}_k)$  of Eq. (25). With active defence (D),  $\tilde{z}_k$  increases until scaling becomes effective, then settles, limiting attacker-realizable accelerations.

The task space comprises planar position  $(x, y)$  and one orientation  $\theta_z$ . The noise covariances are  $\mathbf{R} = 10^{-4} \mathbf{I}$  and  $\mathbf{Q} = \text{blkdiag}_j(q_c \mathbf{Q}_{\text{base}})$ , where  $\mathbf{Q}_{\text{base}} = [T_s^3/3, T_s^2/2; T_s^2/2, T_s]$  and  $q_c = 10^{-2} \text{ rad}^2/\text{s}^3$ , determined with the tuning-knob process.

The nominal task is to maintain the fixed end-effector pose with zero reference velocity and acceleration

$$\bar{\mathbf{p}}_0 = \begin{bmatrix} 1.48 \\ 1.81 \end{bmatrix} \text{ m}, \quad \bar{\mathbf{R}}_0 \text{ from } \mathbf{q}_0 = [0, \frac{\pi}{8}, \dots, \frac{\pi}{8}]^\top.$$

The system is discretized with sampling time  $T_s = 3 \cdot 10^{-3} \text{ s}$ . The passive defence ( $\chi^2$ ) uses a sliding window of  $W = 100$  samples, and it is calibrated for an Average Run Length (ARL) of  $10^6$  samples (i.e., one false positive every 0.95 years under  $\mathcal{H}_0$ ), which is realized through Eq. (13) by setting  $\tau = 779.28$ .

The active defence uses  $f(\cdot)$  of Eq. (25) with  $\beta = (1 - 10^{-3})$ ,  $\psi = (1 - 10^{-6})$ , and  $\gamma = 8$ , determining  $z_p = 50.83$ . For Corollary 2, we have formal guarantees

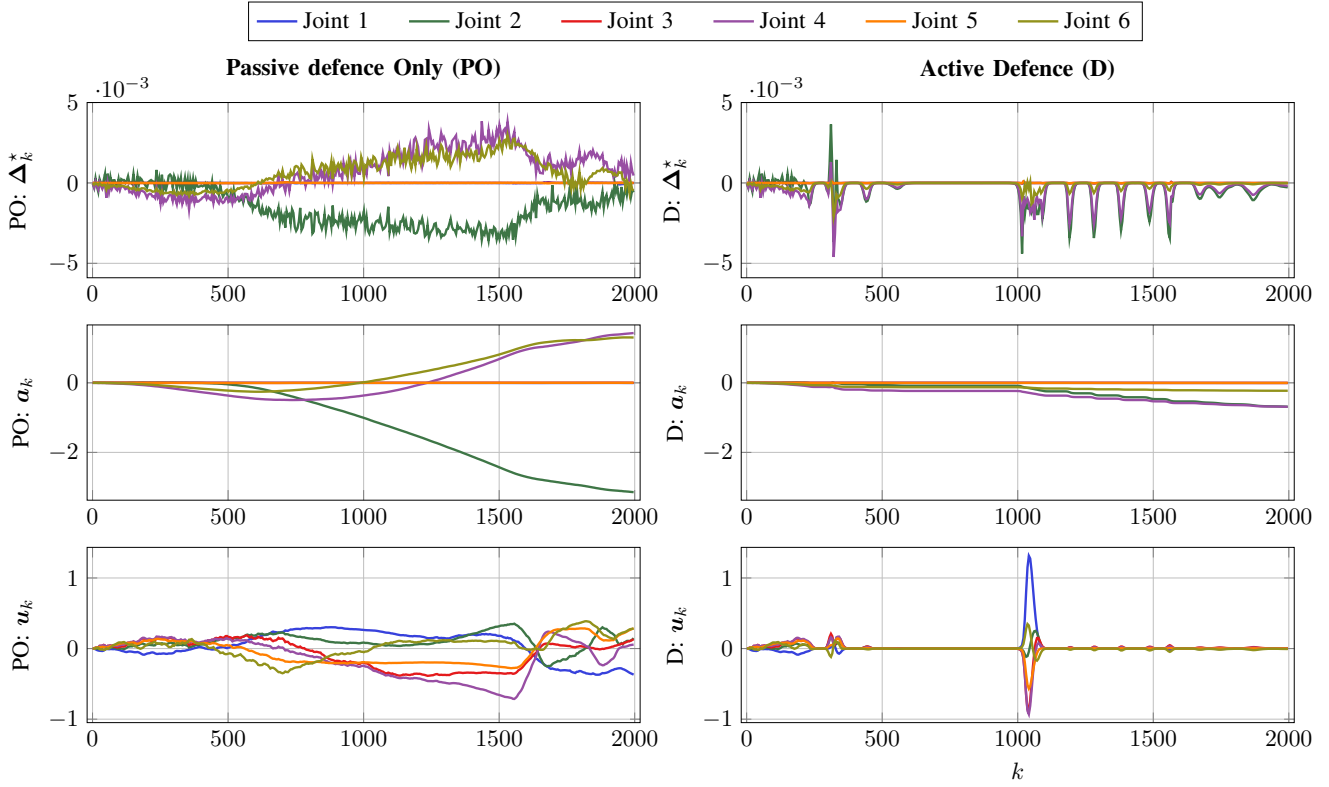


Fig. 4. Adversary and control signals. Optimal increment  $\Delta_k^*$  of Eq. (41), attack sequence  $a_k$  of Eq. (32), and resulting actuation  $u_k$  of Eqs. (9) and (27). Under the proposed active defence (D), the attacker escalates  $\Delta_k^*$  but scaling suppresses the realized  $u_k$ , constraining task-space manipulation.

that, on average and under  $\mathcal{H}_0$ , a single sample of the control command has magnitude below 99.9% of the nominal one every 0.83 hours, which we consider negligible.

The defence adds negligible computational overhead. The QCQP for the attacker (not part of the defence) solves in 40 ms/step with a standard solver (Gurobi).

*Attacker.*: The attacker's goal is to displace the end-effector from  $\bar{\mathbf{p}}_0^A = \bar{\mathbf{p}}_0$  to  $\bar{\mathbf{p}}_{(T-1)}^A = [-2, 1]^T$ . This maneuver is interpolated by quintic polynomials with zero boundary velocities and accelerations, generating the reference trajectory  $\{\bar{\mathbf{p}}_k^A, \dot{\bar{\mathbf{p}}}_k^A, \ddot{\bar{\mathbf{p}}}_k^A\}$  for  $k \in [0, T-1]$ . The anomaly budget is distributed uniformly across time, resulting in a per-step constraint  $\tau' = \tau/T = 7.79$ .

*Evaluation Metrics.*: Let  $\{u_k\}, k = [0, T-1]$  denote a sequence of control commands. The following metric quantifies the control effort:

$$\text{mean}(\{u_k\}) := \frac{1}{T} \sum_k \|u_k\|_2.$$

Let  $\{\bar{\mathbf{p}}_k\}, k = [0, T-1]$  denote a reference hand position trajectory, and  $\{\mathbf{p}_k\}$  the realized one. The following metrics quantify a task deviation:

$$\text{devmax}(\{\bar{\mathbf{p}}_k\}, \{\mathbf{p}_k\}) := \max_k \|\bar{\mathbf{p}}_k - \mathbf{p}_k\|_2,$$

$$\text{devRMS}(\{\bar{\mathbf{p}}_k\}, \{\mathbf{p}_k\}) := \left( \frac{1}{T} \sum_k \|\bar{\mathbf{p}}_k - \mathbf{p}_k\|_2^2 \right)^{1/2}.$$

## B. Main Results

*Undefended (U)*: When no active or passive defence is in place, the end-effector closely follows the malicious trajectory with a small maximum and RMS deviation from the reference (see first column of Table I).

*Passive defence Only (PO)*: In the presence of the  $\chi^2$  ADS using anomaly measure  $z_k$  as in Eq. (11), the attacker needs to satisfy the constraint of Eq. (40). The resulting task deviation mildly increases compared to the Undefended case, and the attacker can still reach its final position goal (see the upper diagrams of Fig. 2 and Fig. 3, and the second column of Table I). Throughout, the anomaly measure  $z_k$  remains below  $\tau'$ , so the passive ADS never fires alarms (see second diagram of Fig. 3).

*With the proposed active defence (D)*: With the command scaling of Eq. (27) in place, in the initial phase ( $k < 250$ ) the optimal attack increment  $\Delta_k^*$  and the deviation  $\|\bar{\mathbf{p}}_k^A - \mathbf{p}_k\|_2$  mirror PO and U, as  $\tilde{z}_k$  is small and  $f(\cdot) \approx 1$  (see Figs. 3 and 4). Subsequently,  $\tilde{z}_k$  increases determining an increasing scaling of command  $u_k$  (see Fig. 4). This incentivizes the attacker to increase injections;  $z_k$  quickly climbs to the per-step limit  $\tau'$  rendering the constraint of Eq. (40) active for a short time.

As the attack progresses, gain reduction  $f(\tilde{z}_k)$  disincentivizes the attacker policy of Section IV-C, since increased anomaly would further strengthen scaling. Consequently,  $\tilde{z}_k$  tends to settle to an asymptotic value at which the acceleration commands are greatly diminished in magnitude.

TABLE I  
TASK DEVIATION AND CONTROL EFFORT METRICS

Metric	U (no defence)	PO ( $\chi^2$ only)	D (active defence)
$\text{devmax}(\{\bar{\mathbf{p}}^A\}, \{\mathbf{p}\})$	0.05	0.28	2.21
$\text{devmax}(\{\bar{\mathbf{p}}\}, \{\mathbf{p}\})$	3.58	3.68	1.5
$\text{devRMS}(\{\bar{\mathbf{p}}^A\}, \{\mathbf{p}\})$	0.02	0.13	1.42
$\text{devRMS}(\{\bar{\mathbf{p}}\}, \{\mathbf{p}\})$	2.22	2.2	0.82
$\text{mean}(\{\mathbf{u}_k\})$	0.44	0.46	0.06

Relative to PO, the realized accelerations are significantly attenuated, yielding a much smaller mean control effort (see Table I).

### C. Comparative Analysis and Discussion

Without defence, or with passive defence only, stealthy FDIAs can precisely steer the end-effector while remaining within the ADS budget; the passive detector is therefore ineffective in isolation.

With the proposed active defence, the attacker faces the following trade-off: increasing injection raises the proposed anomaly-aware command scaling, while reducing injection makes the stealthy attack less effective. Overall, the attack decreases in effectiveness and the proposed anomaly-aware score  $\tilde{z}_k$  settles to an asymptotic value. This plateau represents the sub-optimal injection magnitude the attacker is forced to commit to, akin to an equilibrium in a Stackelberg game. Fig. 3 explains the closed-loop interaction: (i) the attacker increases  $\Delta_k^*$  to counter scaling (Fig. 4, top), (ii) this raises  $z_k$  until the per-step constraint becomes active, (iii) meanwhile  $\tilde{z}_k$  (based on the actuation-projected predictor) tracks the accumulating state discrepancy and drives  $f(\tilde{z}_k)$  down, (iv) the realized  $\mathbf{u}_k$  is therefore attenuated, limiting achievable task-space acceleration and displacement.

An important consideration regards the safety of the proposed active defence in the absence of attacks. By selecting the design confidence  $\psi$  and gain floor  $\beta$  via Corollary 2, unwarranted scaling under  $\mathcal{H}_0$  can be made arbitrarily rare, although statistical fluctuations can in principle still trigger noticeable attenuation. However, because the gain scaling acts on the command  $\mathbf{u}_k$  prior to the inverse-dynamics compensation, it does not alter gravity or Coriolis compensation. Finally, extending the stability analysis to the time-varying case of  $f(\tilde{z}_k)$  is left for future work; promising directions include (i) low-pass filtering  $f(\tilde{z}_k)$ , (ii) enforcing a probabilistic lower bound  $f \geq \beta$ , and (iii) invoking slowly-varying stability results.

## VI. CONCLUSION

This paper addressed the resilience of robotic manipulators against stealthy false data injection attacks (FDIAs). We showed that feedback linearization induces an *integrator vulnerability* that allows persistent sensor corruption to remain undetected by  $\chi^2$  anomaly detectors while driving the end-effector off-task. To counter this threat, we introduced *anomaly-aware command scaling*, a lightweight modification of the control law that attenuates inputs as a function

of an anomaly score derived from a measurement-free, actuation-projected predictor. Our analysis established two key guarantees: (i) probabilistic bounds on actuation loss in nominal operation, enabling minimally invasive deployment, and (ii) preservation of closed-loop stability under bounded attenuation.

On the adversary side, we derived a convex QCQP formulation of the one-step optimal stealthy attack that incorporates the defender's policy, yielding a Stackelberg-type game-theoretic benchmark against which to evaluate defences. Simulation results on a 6-DOF manipulator confirmed that the proposed defence substantially reduces attacker-induced task-space deviations while maintaining nominal tracking performance in the absence of attacks.

Overall, the proposed approach illustrates how active modification of the control law can transform anomaly detection into a practical resilience mechanism for cybersecure robotic manipulation.

## REFERENCES

- [1] A. Humayed, J. Lin, F. Li, and B. Luo, "Cyber-physical systems security—A Survey," *IEEE Internet Things J.*, vol. 4, no. 6, pp. 1802–1831, 2017.
- [2] H. Sandberg, V. Gupta, and K. H. Johansson, "Secure networked control systems," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 5, no. 1, pp. 445–464, 2022.
- [3] Z. Guo, D. Shi, K. H. Johansson, and L. Shi, "Optimal linear cyber-attack on remote state estimation," *IEEE Trans. Control Netw. Syst.*, vol. 4, no. 1, pp. 4–13, 2017.
- [4] A. Intriago, F. Liberati, N. D. Hatzigiorgiou, and C. Konstantinou, "Residual-based detection of attacks in cyber-physical inverter-based microgrids," *IEEE Trans. Power Syst.*, vol. 39, no. 2, pp. 4020–4038, 2024.
- [5] Y. Mo and B. Sinopoli, "Secure control against replay attacks," in *Allerton Conf. Comm., Control, and Comp. (Allerton)*, 2009, pp. 911–918.
- [6] H. Fawzi, P. Tabuada, and S. Diggavi, "Secure estimation and control for cyber-physical systems under adversarial attacks," *IEEE Trans. Autom. Control*, vol. 59, no. 6, pp. 1454–1467, 2014.
- [7] J. Ueda and J. Blevins, "Affine transformation-based perfectly undetectable false data injection attacks on remote manipulator kinematic control with attack detector," *IEEE Robot. Autom. Lett.*, vol. 9, no. 10, pp. 8690–8697, 2024.
- [8] F. E. Tosun, A. M. Teixeira, J. Dong, A. Ahlén, and S. Dey, "Kullback-leibler divergence-based observer design against sensor bias injection attacks in single-output systems," *IEEE Transactions on Information Forensics and Security*, 2025.
- [9] D. Ding, Q.-L. Han, Y. Xiang, X. Ge, and X.-M. Zhang, "A survey on security control and attack detection for industrial cyber-physical systems," *Neurocomputing*, vol. 275, pp. 1674–1683, 2018.
- [10] T. R. C. Murguia, and J. Ruths, "Tuning windowed chi-squared detectors for sensor attacks," in *American Control Conference (ACC)*, 2018, pp. 1752–1757.
- [11] C. Murguia and J. Ruths, "CUSUM and chi-squared attack detection of compromised sensors," in *IEEE Conf. Control Appl. (CCA)*, 2016, pp. 474–480.
- [12] R. M. Murray, Z. Li, and S. Sastry, *A mathematical introduction to robotic manipulation*, 1st ed. Boca Raton: CRC Press, 1994.
- [13] B. Siciliano, L. Sciacivco, L. Villani, and G. Oriolo, *Robotics: Modelling, Planning and Control*, ser. Advanced Textbooks in Control and Signal Processing. Springer London, 2009.
- [14] G. Gallego, C. Cuevas, R. Mohedano, and N. Garcia, "On the Mahalanobis distance classification criterion for multidimensional normal distributions," *IEEE Trans. Signal Process.*, vol. 61, no. 17, pp. 4387–4396, 2013.
- [15] K. M. Ramachandran and C. P. Tsokos, *Sampling distributions*, 3rd ed. Elsevier, 2021, pp. 147–177.
- [16] G. Gualandi, "Manipulatorsgainscaling," 2026, v1.0.1. [Online]. Available: <https://doi.org/10.5281/zenodo.18877674>