

Active Defense Against False Data Injection Attacks in Robotic Manipulators

Gabriele Gualandi* Carl Mikael Larsson*
Alessandro V. Papadopoulos*

* *Mälardalen University, Västerås, Sweden*
(e-mail: gabriele.gualandi@mdu.se).

Abstract: Robotic systems are vulnerable to *False Data Injection Attacks* (FDIAs), where adversaries corrupt sensor signals to gain malicious control. Feedback linearization exposes robotic systems to *integrator vulnerability*, exposing them to *stealthy* attacks that can cause significant deviations in end-effector behavior without raising alarms. This paper addresses the resilience of manipulators against finite-horizon FDIAs by formalizing two defense methods, namely *anomaly-aware virtual damping* and *manipulability reduction*, with probabilistic guarantees on nominal task execution. Simulations on a 7-DOF redundant manipulator show that the proposed defense substantially reduces the impact of FDIA compared to using solely a threshold-based ADS like the χ^2 , while preserving nominal task performance in the absence of attack.

1. 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 (Humayed et al., 2017; Sandberg et al., 2022). 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 detection while still driving the end-effector away from its intended task (Guo et al., 2017; Intriago et al., 2024).

A well-studied class of integrity attacks are *False Data Injection Attacks* (FDIAs) targeting sensing, 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 (Mo and Sinopoli, 2009; Ueda and Blevins, 2024; Fawzi et al., 2014) 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 (Ueda and Blevins, 2024). One perspective has been largely overlooked which relates to a key robotics-specific vulnerability: the widespread use of *feedback linearization* reduces manipulator dynamics to double integrators, which in turn introduces an *integrator vulnerability* (Tosun et al., 2025). As a consequence, persistent sensor corruption can silently accumulate in the closed-loop system.

This paper addresses this gap by introducing two active defense strategies that limit the impact of finite-horizon FDIA targeting sensors for feedback-linearized robotic arms. The proposed defense methods leverages an *actuation-projected anomaly score* that ignores sensing and hence is immune to sensor corruption. Said anomaly score is used in two ways. Firstly, we introduce *virtual damping* for the joints, to dissipate kinetic energy as anomalies grow, thereby reducing the adversary’s ability to steer the manipulator. Secondly, we reconfigure the manipulator’s null space to *reduce manipulability* in the estimated attack hand direction, which increases the joint-level anomaly required for a given hand-level displacement. Since the strength of the virtual damping, and the detection of a standard χ^2 ADS are both driven by joint-level anomaly, null-space reconfiguration can further limit the impact of FDIA.

The contributions of this work are the following:

- (1) We formalize stealthy FDIAs against feedback-linearized manipulators, exposing their integrator vulnerability, and show that the attacker’s one-step optimal strategy reduces to a convex QCQP.
- (2) We propose *anomaly-aware virtual damping*, which attenuates joint velocities based on a measurement-free actuation-projected state predictor, with probabilistic guarantees on bounded attenuation in nominal operation. Furthermore, we propose to *decrease manipulability* in the estimated attack direction to further reduce the impact of FDIA, with no impact on task performance in nominal operation.
- (3) Simulations on a 7-DOF real-world robot shows that the defenses significantly limits FDIA compared to threshold-based ADSs like the χ^2 .

1.1 Related Work

Threshold-based detectors, particularly χ^2 tests on Kalman innovations, are a classical tool for monitoring

* This work was supported by the Knowledge Foundation (KKS).

CPS integrity (Ding et al., 2018). These methods provide statistical guarantees on false-alarm rates and are widely used in industrial practice. Extensions include adaptive thresholds (Tunga et al., 2018) and sequential schemes such as CUSUM tests (Murguia and Ruths, 2016). However, such ADS are inherently *passive*: they may detect anomalies but do not alter the control policy, leaving the system structure unchanged.

Theoretical studies of FDIAs in networked systems have characterized undetectability conditions (Mo and Sinopoli, 2009), affine attack structures (Ueda and Blevins, 2024), and optimal attack strategies (Guo et al., 2017). These works typically assume 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 (Intriago et al., 2024).

Control systems with a LTI plant having integral action are subject to *integrator vulnerability* (Tosun et al., 2025), where constant attack injections get absorbed by any linear observer, resulting in attacks only being detectable during transients. Classical passive threshold-based detection cannot prevent a stealth FDIA from increasingly deceiving a state estimator.

In (Gualandi and Papadopoulos, 2026), we demonstrated that feedback linearization induces integrator vulnerability in robotic systems, and we proposed gain scaling based on a state-projected (sensing-free) anomaly score to limit the accumulation of kinetic energy during a FDIA. While gain scaling mitigates this energy accumulation, it is fundamentally unable to dissipate it. Therefore, this work introduces a novel *virtual damping* strategy designed to actively dissipate the kinetic energy generated by a FDIA, while maintaining probabilistic performance guarantees during nominal operation. Furthermore, we propose directional *manipulability reduction* to compel the adversary to increase its signal injection for a given task-space displacement, thereby further curtailing its control authority. This establishes kinematic redundancy not only as a mechanism for performance enhancement but also as a tool for cyber-physical security. Although redundant manipulators have long leveraged null-space projections for secondary tasks, such as obstacle clearance (Siciliano and Slotine, 1991), their application to adversarial robustness remains largely unexplored. To the best of our knowledge, exploiting redundancy to *reduce manipulability along an adversarial direction* as an active security defense has not been previously investigated.

2. SYSTEM MODEL

2.1 Closed-Loop Architecture

We consider a robotic manipulator operating in closed loop with a state estimator and a task-space controller. At each discrete time step $k \in \mathbb{Z}_{\geq 0}$, the plant output \mathbf{y}_k is corrupted by injected signal \mathbf{a}_k to yield measurement

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

The plant has discrete-time linear dynamics with input saturation

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

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

with state \mathbf{x}_k , process noise $\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$, measurement noise $\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$, and input \mathbf{u}_k saturated by the $\text{sat}(\cdot)$ operator within limits $[\mathbf{u}_{\min}, \mathbf{u}_{\max}]$, $\mathbf{u}_{\min} < \mathbf{0}$, $\mathbf{u}_{\max} > \mathbf{0}$. 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 is the innovation (residual). We assume (\mathbf{A}, \mathbf{C}) is detectable, ensuring the existence of a unique stabilizing solution \mathbf{P} to the DARE

$$\mathbf{P} = \mathbf{A}\mathbf{P}\mathbf{A}^\top + \mathbf{Q} - \mathbf{A}\mathbf{P}\mathbf{C}^\top (\mathbf{C}\mathbf{P}\mathbf{C}^\top + \mathbf{R})^{-1} \mathbf{C}\mathbf{P}\mathbf{A}^\top. \quad (4)$$

Here \mathbf{P} is the steady-state covariance of the signal $\mathbf{e}_k = \mathbf{x}_k - \hat{\mathbf{x}}_k$. The steady-state covariance of \mathbf{r}_k is

$$\mathbf{\Sigma} = \mathbf{C}\mathbf{P}\mathbf{C}^\top + \mathbf{R} > \mathbf{0} \quad (5)$$

(hence invertible), and the steady-state gain is

$$\mathbf{L} = \mathbf{A}\mathbf{P}\mathbf{C}^\top \mathbf{\Sigma}^{-1}. \quad (6)$$

In the remainder of the paper, symbols with the hat $\hat{\cdot}$ are quantities determined using the estimated state $\hat{\mathbf{x}}$, which is subject to FDIA.

2.2 Manipulator Model and Task-Space Control

We consider an n -DOF robotic manipulator with joint variable $\mathbf{q} \in \mathbb{R}^n$, state $\mathbf{x}_k = [\mathbf{q}_k^\top \dot{\mathbf{q}}_k^\top]^\top \in \mathbb{R}^{2n}$ and measurement vector $\mathbf{y}_k \in \mathbb{R}^p$ (assuming full state measurement, $p = 2n$). The continuous-time joint-space dynamics with joint torques $\boldsymbol{\tau}$ is:

$$\mathbf{B}(\mathbf{q})\ddot{\mathbf{q}} + \boldsymbol{\nu}(\mathbf{q}, \dot{\mathbf{q}}) = \boldsymbol{\tau}, \quad (7)$$

where $\mathbf{B}(\mathbf{q}) > \mathbf{0}$ is the inertia matrix and $\boldsymbol{\nu}$ collects Coriolis, centrifugal, and gravity terms. With inverse-dynamics compensation $\boldsymbol{\tau} = \mathbf{B}(\mathbf{q})\mathbf{u} + \boldsymbol{\nu}(\mathbf{q}, \dot{\mathbf{q}})$, the joint dynamics reduce to decoupled double integrators (*virtual masses system*), $\ddot{\mathbf{q}} = \mathbf{u}$, whose discrete-time form is (2a) and whose residual is $\mathbf{r}_k \in \mathbb{R}^p$ from (3b).

The controller operates in task space using twist control (Siciliano et al., 2009), tracking position $\{\bar{\mathbf{p}}_k, \dot{\bar{\mathbf{p}}}_k, \ddot{\bar{\mathbf{p}}}_k\}$ and orientation $\{\bar{\mathbf{R}}_k, \bar{\boldsymbol{\omega}}_k, \dot{\bar{\boldsymbol{\omega}}}_k\}$ references, with $\bar{\mathbf{R}}_k \in SO(3)$, $\bar{\mathbf{p}}_k \in \mathbb{R}^3$ and $\bar{\boldsymbol{\omega}}_k \in \mathbb{R}^3$. A PD+feedforward law provides task accelerations:

$$\mathbf{u}_k^{\text{PD}} = \begin{bmatrix} \ddot{\bar{\mathbf{p}}}_k + K_{pp} \mathbf{e}_{p,k} + K_{dp} \dot{\mathbf{e}}_{p,k} \\ \dot{\bar{\boldsymbol{\omega}}}_k + K_{po} \mathbf{e}_{o,k} + K_{do} \mathbf{e}_{\omega,k} \end{bmatrix}, \quad (8)$$

where: $\mathbf{e}_{p,k} = \bar{\mathbf{p}}_k - \hat{\mathbf{p}}_k$ and $\mathbf{e}_{o,k} = \sin(\hat{\theta}_k/2) \hat{\mathbf{r}}_k$ are respectively the position and orientation errors; $\mathbf{e}_{o,k}$ is computed using the angle-axis pair $(\hat{\theta}_k, \hat{\mathbf{r}}_k)$ from matrix $\bar{\mathbf{R}}_k \bar{\mathbf{R}}_k^\top$, $\hat{\theta}_k \in [0, \pi]$; $\mathbf{e}_{\omega,k} = \bar{\boldsymbol{\omega}}_k - \hat{\boldsymbol{\omega}}_k$; and the gains $K_{pp}, K_{dp}, K_{po}, K_{do}$ are synthesized via discrete-time LQR.

Joint accelerations are computed via pseudoinverse $\mathbf{J}^*(\hat{\mathbf{q}})$ of the geometric Jacobian (Siciliano et al., 2009) as:

$$\mathbf{u}_k^{\text{nom}} = \mathbf{J}_k^* (\mathbf{u}_k^{\text{PD}} - \dot{\mathbf{J}}_k \dot{\hat{\mathbf{q}}}) + \mathbf{u}_k^{\text{sec}} \quad (9)$$

where $\mathbf{u}_k^{\text{sec}}$ is later used for an active defense lying in the null-space of the Jacobian, therefore not interfering with the primary task under \mathcal{H}_0 (no attack). In the absence of active defenses, in (2a) we set $\mathbf{u}_k = \mathbf{J}_k^* (\mathbf{u}_k^{\text{PD}} - \dot{\mathbf{J}}_k \dot{\hat{\mathbf{q}}})$.

2.3 χ^2 Anomaly Detection System

The system incorporates a detector to distinguish between the null hypothesis of normal operation, \mathcal{H}_0 , and the alternative hypothesis of an attack, \mathcal{H}_1 . The χ^2 detector monitors the statistic (Mahalanobis distance):

$$z_k = \mathbf{r}_k^\top \boldsymbol{\Sigma}^{-1} \mathbf{r}_k \quad (10)$$

with $\boldsymbol{\Sigma}$ as in (5) and \mathbf{r}_k as in (3b). During execution, an alarm is triggered if $z_k > \tau$, where $\tau > 0$ is a parameter (*threshold*). Under \mathcal{H}_0 , the residual $\mathbf{r}_k \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$, hence z_k follows a $\chi^2(p)$ distribution, where $p = \dim(\mathbf{r}_k)$ (Ding et al., 2018). Given a desired per-step false-alarm probability $\alpha \in (0, 1)$ (equivalently, a desired mean inter-alarm interval $\text{ARL} = 1/\alpha$), setting

$$\tau = F_{\chi^2(p)}^{-1}(1 - \alpha) = 2P^{-1}\left(1 - \alpha, \frac{p}{2}\right), \quad (11)$$

ensures that $\mathbb{P}(z_k > \tau \mid \mathcal{H}_0) = \alpha$, where $F_{\chi^2(p)}^{-1}$ is the inverse CDF of the χ^2 distribution with p degrees of freedom, and P^{-1} denotes the inverse of the regularized lower incomplete gamma function (Tunga et al., 2018).

2.4 Assumptions

There are the following assumptions, under which both the defense and attack strategies are developed.

Assumption 1. (Adversary knowledge). The attacker has full knowledge of the plant, controller, state estimator, ADS and active defense system(s).

Assumption 2. (Adversary capabilities). The only attack surface is additive injection into sensor measurements FDIA targeting sensors, as in (1).

Assumption 3. (Adversary goal). The attacker's goal is to control the position of the end-effector in task space while remaining undetected by any threshold-based ADS such as the χ^2 (stealth FDIA).

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 has converged (covariances and gains are constant).

3. PROPOSED DEFENCE METHODS

3.1 Attack Estimation

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

$$\tilde{\mathbf{x}}_0 = \hat{\mathbf{x}}_0. \quad (12)$$

For subsequent steps, $\tilde{\mathbf{x}}_k$ follows the open-loop predictor

$$\tilde{\mathbf{x}}_{k+1} = \mathbf{A} \tilde{\mathbf{x}}_k + \mathbf{B} \mathbf{u}_k, \quad k \geq 0. \quad (13)$$

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

$$\tilde{\mathbf{r}}_k = \hat{\mathbf{x}}_k - \tilde{\mathbf{x}}_k \in \mathbb{R}^{2n}. \quad (14)$$

Theorem 1. (Actuation-projected residual covariance).

Under \mathcal{H}_0 (no attack) the covariance of the residual in

(14) after k steps from initialization, $\boldsymbol{\Sigma}_{\tilde{\mathbf{r}},k} \in \mathbb{R}^{2n \times 2n}$, is the (2, 2) block of the matrix $\mathbf{P}_{z,k} \in \mathbb{R}^{4n \times 4n}$, where blocks have size $2n \times 2n$, obtained by iterating

$$\mathbf{P}_{z,j+1} = \mathbf{F} \mathbf{P}_{z,j} \mathbf{F}^\top + \boldsymbol{\Pi}, \quad j = 0, \dots, k-1, \quad (15)$$

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

$$\mathbf{F} = \begin{bmatrix} \mathbf{A} - \mathbf{L}\mathbf{C} & \mathbf{0} \\ \mathbf{L}\mathbf{C} & \mathbf{A} \end{bmatrix} \in \mathbb{R}^{4n \times 4n}, \quad (16)$$

$$\mathbf{Z} = \begin{bmatrix} \mathbf{I} & -\mathbf{L} \\ \mathbf{0} & \mathbf{L} \end{bmatrix} \in \mathbb{R}^{4n \times (2n+p)}, \quad (17)$$

$$\boldsymbol{\Pi} = \mathbf{Z} \text{diag}(\mathbf{Q}, \mathbf{R}) \mathbf{Z}^\top \in \mathbb{R}^{4n \times 4n}. \quad (18)$$

and \mathbf{P} is defined in (4).

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 $\tilde{\mathbf{r}}_k := \hat{\mathbf{x}}_k - \tilde{\mathbf{x}}_k$ are:

$$\begin{aligned} \mathbf{e}_{k+1} &= (\mathbf{A} - \mathbf{L}\mathbf{C})\mathbf{e}_k + \mathbf{w}_k - \mathbf{L}\mathbf{v}_k, \\ \tilde{\mathbf{r}}_{k+1} &= \mathbf{A}\tilde{\mathbf{r}}_k + \mathbf{L}\mathbf{C}\mathbf{e}_k + \mathbf{L}\mathbf{v}_k. \end{aligned}$$

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

$$\mathbf{z}_{k+1} = \mathbf{F}\mathbf{z}_k + \mathbf{Z}\boldsymbol{\eta}_k,$$

with \mathbf{F}, \mathbf{Z} as in (16), (17) and $\boldsymbol{\eta}_k \triangleq [\mathbf{w}_k^\top, \mathbf{v}_k^\top]^\top$. The covariance propagates as

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

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

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

with $\boldsymbol{\Pi}$ as in (18). At synchronization $k = 0$, for (12) and Assumption 4, we have $\tilde{\mathbf{r}}_0 = \mathbf{0}$ and $\text{Cov}(\mathbf{e}_0) = \mathbf{P}$, hence

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

Iterating (15) k times yields $\mathbf{P}_{z,k}$; the desired covariance $\boldsymbol{\Sigma}_{\tilde{\mathbf{r}},k} \in \mathbb{R}^{2n \times 2n}$ is its lower-right $2n \times 2n$ submatrix. \square

Theorem 2. (Confidence for the Anomaly Measure).

Consider the statistic (*projected anomaly score*)

$$\tilde{z}_k = \tilde{\mathbf{r}}_k^\top \boldsymbol{\Sigma}_{\tilde{\mathbf{r}},k}^{-1} \tilde{\mathbf{r}}_k. \quad (19)$$

Choosing $z_x = F_{\chi^2(2n)}^{-1}(\psi)$ ensures $\mathbb{P}(\tilde{z}_k \leq z_x \mid \mathcal{H}_0) = \psi$ for any given $z_x > 0$ and probability ψ .

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

3.2 Virtual Damping Active Defense

We propose virtual damping to enforce a passivity margin on the virtual masses system based on \tilde{z}_k as in (19).

Let $F > 0$ be the desired asymptotic power dissipation, $\gamma > 0$ a design parameter, and let $z_x > 0$ be a design abscissa. Let $\beta \in (0, F)$ be the desired power dissipation ratio at $\tilde{z} = z_x$. We define the *target power dissipation*

ratio, $g(\tilde{z})$, as a smooth, monotonically increasing sigmoid function:

$$g(\tilde{z}_k) = F \left(1 - \exp \left[- \left(\frac{\tilde{z}_k}{z_s} \right)^\gamma \right] \right), \quad (20)$$

$$z_s = z_x \left(- \ln \left(1 - \frac{\beta}{F} \right) \right)^{-1/\gamma}, \quad (21)$$

satisfying $g(0) = 0$, $g(\tilde{z}_k) \geq 0$, $\lim_{z \rightarrow \infty} g(z) = F$ and $g(z_x) = \beta$.

Assumption 5. (Model fidelity). We assume that $\tilde{\mathbf{x}}$ always remains accurate over a finite-horizon FDIA on sensors i.e.,

$$\tilde{\mathbf{x}}_k \approx \mathbf{x}_k, k \in (0, T). \quad (22)$$

Assumption 5 is only needed for the accuracy of the defense control laws based on actuation-projected state $\tilde{\mathbf{x}}_k$, as developed in the rest of this work. It is reasonable when process noise remains small over the attack horizon T (which is limited by mechanisms such as periodic software rejuvenation, key rotation, re-establishment of secure channels, or moving-target defenses) and when $\tilde{\mathbf{x}}_k$ is regularly resynchronized with $\hat{\mathbf{x}}_k$ as in (12).

Theorem 3. (Adaptive damping). Consider a robotic manipulator with nominal joint-space control input $\mathbf{u}_k^{\text{nom}}$ (9) subject to the saturation operator $\text{sat}(\cdot)$ (2a). Let the (projected) nominal power of the j -th joint be defined as

$$\tilde{P}_{k,j}^{\text{nom}} = \text{sat}(\mathbf{u}_k^{\text{nom}})_j \cdot \dot{\tilde{q}}_{k,j} \quad (23)$$

where $\tilde{\mathbf{q}}_k$ is extracted from $\tilde{\mathbf{x}}_k$. Let $\tilde{z}_k \in [0, \infty)$ denote the anomaly score defined in (19), and let $g : [0, \infty) \rightarrow [0, \infty)$ be a non-decreasing function satisfying $g(0) = 0$.

Define the ideal, unconstrained adaptive damping command, $\mathbf{u}_k^{\text{d,ideal}}$, on a per-joint basis as

$$\mathbf{u}_{k,j}^{\text{d,ideal}} = \begin{cases} -g(\tilde{z}_k) \frac{|\tilde{P}_{k,j}^{\text{nom}}|}{\dot{\tilde{q}}_{k,j}}, & \text{if } |\dot{\tilde{q}}_{k,j}| > \epsilon, \\ 0, & \text{otherwise,} \end{cases} \quad (24)$$

with $g(\cdot)$ an arbitrary non-decreasing function satisfying $g(0) = 0$ (such as the one defined in (20)) and $\epsilon \approx 0$.

Let us define the interval $[\mathbf{h}_k^-, \mathbf{h}_k^+]$, where

$$\mathbf{h}_k^- = \mathbf{u}_{\min} - \text{sat}(\mathbf{u}_k^{\text{nom}}) \in [\mathbf{u}_{\min} - \mathbf{u}_{\max}, \mathbf{0}], \quad (25a)$$

$$\mathbf{h}_k^+ = \mathbf{u}_{\max} - \text{sat}(\mathbf{u}_k^{\text{nom}}) \in [\mathbf{0}, \mathbf{u}_{\max} - \mathbf{u}_{\min}], \quad (25b)$$

are the actuator headrooms available for the defense.

Define the *maximal non-saturating damping* command, \mathbf{u}_k^{d} , as:

$$\mathbf{u}_k^{\text{d}} = \max(\mathbf{h}_k^-, \min(\mathbf{h}_k^+, \mathbf{u}_k^{\text{d,ideal}})). \quad (26)$$

If the final control law is

$$\mathbf{u}_k = \text{sat}(\mathbf{u}_k^{\text{nom}}) + \mathbf{u}_k^{\text{d}}, \quad (27)$$

then, under Assumption 5:

Claim 1. (Stability). Command \mathbf{u}_k^{d} acting on the virtual masses system has non-positive joint-wise power, i.e.:

$$\mathbf{u}_{k,j}^{\text{d}} \cdot \dot{q}_{k,j} \leq 0, \quad \forall j = 1, \dots, n.$$

Hence, for the fundamental results in passivity theory (Van der Schaft, 2000), it does not decrease the stability margin of the closed loop system.

Claim 2. (Probabilistic Magnitude Bound). Let $\tilde{z}_k \sim \chi^2(2n)$ under \mathcal{H}_0 , and let $z_x = F_{\chi^2(2n)}^{-1}(\psi)$ denote the inverse CDF at confidence level $\psi \in (0, 1)$. If design

parameters of $g(\cdot)$ are chosen such that $g(z_x) \leq \beta$ for some $\beta > 0$, then with probability at least ψ , the magnitude of the total power dissipated by the defense is bounded by a fraction β of the total absolute nominal power:

$$\sum_{j=1}^n |\mathbf{u}_{k,j}^{\text{d}} \cdot \dot{\tilde{q}}_{k,j}| \leq \beta \sum_{j=1}^n |\tilde{P}_{k,j}^{\text{nom}}|. \quad (28)$$

Proof. For Assumption 5, any projected power (computed using $\dot{\tilde{\mathbf{q}}}$) coincides with its real version (from $\dot{\mathbf{q}}$).

To prove Claim 1, let us denote the per-joint power of the ideal command as $\tilde{P}_{k,j}^{\text{d,ideal}} = \mathbf{u}_{k,j}^{\text{d,ideal}} \cdot \dot{\tilde{q}}_{k,j}$. By substituting (24) into this definition, and observing that $\dot{\tilde{q}}_{k,j}/\dot{q}_{k,j} = 1$, we obtain:

$$\tilde{P}_{k,j}^{\text{d,ideal}} = -g(\tilde{z}_k) |\tilde{P}_{k,j}^{\text{nom}}| \leq 0, \quad (29)$$

since $g(\cdot) \geq 0$ and $|\tilde{P}_{k,j}^{\text{nom}}| \geq 0$. This proves that the ideal command is per-joint passive. Next, consider the final command $\mathbf{u}_{k,j}^{\text{d}}$ from (26). Since $\text{sat}(\mathbf{u}_k^{\text{nom}})$ is within actuator limits $[\mathbf{u}_{\min}, \mathbf{u}_{\max}]$, the headroom intervals $[\mathbf{h}_k^-, \mathbf{h}_k^+]$ from (25a)–(25b) are guaranteed to contain the origin. Therefore, the clipping operation in (26) cannot change the sign, nor increase the absolute value, of $\mathbf{u}_{k,j}^{\text{d,ideal}}$, leading to:

$$|\mathbf{u}_{k,j}^{\text{d}} \cdot \dot{\tilde{q}}_{k,j}| \leq |\mathbf{u}_{k,j}^{\text{d,ideal}} \cdot \dot{\tilde{q}}_{k,j}|, \quad (30)$$

$$\text{sign}(\mathbf{u}_{k,j}^{\text{d}} \cdot \dot{\tilde{q}}_{k,j}) = \text{sign}(\mathbf{u}_{k,j}^{\text{d,ideal}} \cdot \dot{\tilde{q}}_{k,j}). \quad (31)$$

By (30) and (31), the passivity condition (29) still holds after clipping, hence (26) is per-joint passive.

To prove Claim 2, consider the absolute total power of the ideal defense command. Substituting (29) into the definition of total power yields:

$$\sum_{j=1}^n |\mathbf{u}_{k,j}^{\text{d,ideal}} \cdot \dot{\tilde{q}}_{k,j}| = \sum_{j=1}^n |-g(\tilde{z}_k) |\tilde{P}_{k,j}^{\text{nom}}|| = g(\tilde{z}_k) \sum_{j=1}^n |\tilde{P}_{k,j}^{\text{nom}}|. \quad (32)$$

Summing (30) over $j = 1, \dots, n$ and substituting the total ideal power from (32), we obtain the total power bound:

$$\sum_{j=1}^n |\mathbf{u}_{k,j}^{\text{d}} \cdot \dot{\tilde{q}}_{k,j}| \leq \sum_{j=1}^n |\mathbf{u}_{k,j}^{\text{d,ideal}} \cdot \dot{\tilde{q}}_{k,j}| = g(\tilde{z}_k) \sum_{j=1}^n |\tilde{P}_{k,j}^{\text{nom}}|. \quad (33)$$

Under \mathcal{H}_0 , the anomaly score follows $\tilde{z}_k \sim \chi^2(2n)$. By the definition of the quantile $z_x = F_{\chi^2(2n)}^{-1}(\psi)$, we have $\mathbb{P}(\tilde{z}_k \leq z_x) = \psi$. On this event, since $g(\cdot)$ is non-decreasing, we have $g(\tilde{z}_k) \leq g(z_x)$. By design, $g(z_x) \leq \beta$. Substituting this into (33) yields

$$\sum_{j=1}^n |\mathbf{u}_{k,j}^{\text{d}} \cdot \dot{\tilde{q}}_{k,j}| \leq \beta \sum_{j=1}^n |\tilde{P}_{k,j}^{\text{nom}}|.$$

with probability ψ . Therefore, with probability ψ , the total dissipated power is bounded by a fraction β of the nominal absolute power, which completes the proof. \square

Claim 2 of Thm. 3 provides a formal method for tuning the defense's nominal impact. By choosing $z_x = F_{\chi^2(2n)}^{-1}(\psi)$, one can guarantee with an arbitrarily high probability ψ (e.g., $\psi = 0.99$) that the power of the adaptive friction will be less than a desired small fraction β (e.g., $\beta = 0.01$)

of the nominal control power, so that the defense can be configured to be minimally invasive under \mathcal{H}_0 . A value of $F > 1$ enforces asymptotic strict passivity of the virtual masses system as the projected anomaly \tilde{z}_k grows.

3.3 Manipulability Reduction Active Defense

We now present a defense strategy that leverages the manipulator's kinematic redundancy. Specifically, we reduce manipulability in the estimated attack direction to increase the required anomaly injected for a given hand displacement. We restrict attack direction estimation to the translational DOFs, as these are most relevant to adversarial manipulation. The estimated attack direction is defined as

$$\mathbf{d}_k = \frac{\tilde{\mathbf{p}}_k - \bar{\mathbf{p}}_k + \epsilon}{\|\tilde{\mathbf{p}}_k - \bar{\mathbf{p}}_k + \epsilon\|_2}, \quad (34)$$

where $\mathbf{d}_k \in \mathbb{R}^3$, $\|\mathbf{d}_k\|_2 = 1$, $\tilde{\mathbf{p}}_k$ is the hand position from $\hat{\mathbf{q}}_k$, and $\epsilon \approx 0$ avoids degeneracy.

The positional manipulability measure in direction \mathbf{d}_k is

$$\eta_k = \mathbf{d}_k^\top \mathbf{M}_k \mathbf{d}_k, \quad (35)$$

where $\mathbf{M}_k = \tilde{\mathbf{J}}_{\mathbf{p},k} \tilde{\mathbf{J}}_{\mathbf{p},k}^\top$ is the directional manipulability matrix, and $\tilde{\mathbf{J}}_{\mathbf{p},k}$ is the positional geometric Jacobian from the projected joint configuration $\hat{\mathbf{q}}_k$.

We adopt the redundancy resolution scheme based on projected gradient method (De Luca and Oriolo, 1991) with the goal of reducing manipulability along \mathbf{d}_k , using cost function:

$$\mathcal{C}_k = \frac{1}{2} \eta_k^2. \quad (36)$$

Treating \mathbf{d}_k as quasi-static allows for the gradient of the cost to be approximated as

$$\nabla_{\mathbf{q}} \mathcal{C}_k \approx \eta_k \begin{bmatrix} \mathbf{d}_k^\top \frac{\partial \mathbf{M}_k}{\partial \mathbf{q}^{(1)}} \mathbf{d}_k \\ \vdots \\ \mathbf{d}_k^\top \frac{\partial \mathbf{M}_k}{\partial \mathbf{q}^{(n)}} \mathbf{d}_k \end{bmatrix}. \quad (37)$$

We introduce the following secondary null-space command for control law $\mathbf{u}_k^{\text{nom}}$ (9):

$$\mathbf{u}_k^{\text{sec}} = \begin{cases} \mathbf{0}, & \text{if } \nabla_{\mathbf{q}} \mathcal{C}_k \approx 0 \\ -\nu_k (\mathbf{I} - \hat{\mathbf{J}}_k^\dagger \hat{\mathbf{J}}_k) \nabla_{\mathbf{q}} \mathcal{C}_k, & \text{otherwise} \end{cases} \quad (38)$$

where $\hat{\mathbf{J}}_k^\dagger$ is the Moore–Penrose inverse, $\nu_k \in (0, \nu_{\max}]$ is a step size selected by an Armijo line search.

To penalize individual per-joint motion and maintain low manipulability when \mathcal{C} is stationary i.e., $\nabla_{\mathbf{q}} \mathcal{C}_k \approx 0$, we use the following weighted pseudoinverse for $\mathbf{u}_k^{\text{nom}}$ (9):

$$\mathbf{J}^* = \mathbf{W}_k^{-1} \hat{\mathbf{J}}_k^\top (\hat{\mathbf{J}}_k \mathbf{W}_k^{-1} \hat{\mathbf{J}}_k^\top)^{-1}, \quad (39)$$

with weighting

$$\mathbf{W}_k = \begin{cases} \mathbf{D} + \alpha (\nabla_{\mathbf{q}}^2 \mathcal{C}_k + \mu_k \mathbf{I}), & \text{if } \nabla_{\mathbf{q}} \mathcal{C}_k \approx 0 \\ \mathbf{D}, & \text{otherwise} \end{cases} \quad (40)$$

where $\mathbf{D} \succ 0$ is diagonal, $\alpha \geq 0$ is a design scalar and $\mu_k \geq 0$ is chosen opportunely to ensure $\mathbf{W}_k \succ 0$.

Theorem 4. (Existence of \mathbf{J}^* and C^1 -continuity of \mathbf{W}_k).
Choosing:

$$\mu_k = \log(1 + \exp(\epsilon - \lambda_{\min})), \quad (41)$$

for (40), where $\epsilon > 0$, $\epsilon \approx 0$, and λ_{\min} is the minimum eigenvalue of $\nabla_{\mathbf{q}}^2 \mathcal{C}_k$; guarantees existence and uniqueness of \mathbf{J}^* , and C^1 -continuity of \mathbf{W}_k .

Proof. The Hessian $\nabla_{\mathbf{q}}^2 \mathcal{C}$ is symmetric by Schwarz's theorem. Adding $\max(0, \epsilon - \lambda_{\min}) \mathbf{I}$ to $\nabla_{\mathbf{q}}^2 \mathcal{C}$ shifts all eigenvalues to be positive (Liao and Shoemaker, 1991), ensuring $\mathbf{W}_k \succ 0$, which guarantees the existence and uniqueness of \mathbf{J}^* . Subjecting the above max operation to a standard softplus function provides (41). Being the eigenvalues of \mathbf{W}_k continuous in \mathbf{q} , \mathbf{W}_k is C^1 -continuous in \mathbf{q} by construction of (40). \square

Choosing a large α makes the term $(\nabla_{\mathbf{q}}^2 \mathcal{C}_k + \mu_k \mathbf{I})$ dominate, while increasing diagonal entries of \mathbf{D} penalizes per-joint motion. To enforce the physical limits $[\mathbf{u}_{\min}, \mathbf{u}_{\max}]$ through the saturation operator $\text{sat}(\cdot)$ (2a) without distorting the commanded motion direction, \mathbf{u}_k is computed from $\mathbf{u}_k^{\text{nom}}$ (9) via hierarchical magnitude-based scaling. The primary task is prioritized by reserving a fixed quota (e.g., 30%) of the actuation capability for the secondary task $\mathbf{u}_k^{\text{sec}}$.

4. OPTIMAL STEALTH ATTACK

Under the assumptions presented in Section 2.4, we formalize the attacker's strategy to control the end-effector while remaining stealth against threshold-based ADS detection. We envisage an attacker that comprehensively simulates the dynamics of all involved system in the absence of noise, like the plant, controller and proposed active defense system(s).

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

Let $\mathcal{Z}_{k,j}(\mathbf{a})$ denote the j -step-ahead prediction of the closed-loop system, starting from time k , under the attack sequence $[\mathbf{a}, \mathbf{0}, \dots, \mathbf{0}]$ (the only non-zero injection occurs at time k). For a generic vector \mathbf{v} , we define

$$\mathbf{v}_{k+j}^{\text{SIM}} := \pi_{\mathbf{v}}(\mathcal{Z}_{k,j}(\mathbf{a})), \quad (42)$$

where $\pi_{\mathbf{v}}(\cdot)$ extracts the vector \mathbf{v} at time $k+j$ from the simulated evolution.

Under the closed-loop dynamics in (2a), (2b), (3a), and (3b), an injected signal at time k , \mathbf{a}_k , influences the end-effector acceleration *two* time steps later. Specifically,

$$\ddot{\mathbf{p}}_{k+2}^{\text{SIM}} = \pi_{\ddot{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_k)). \quad (43)$$

At each time step, the attacker's high-level objective is to make this predicted acceleration match a desired target acceleration $\ddot{\mathbf{p}}_{k+2}^{\text{A}}$. This is realized by minimizing

$$\frac{1}{2} \|\ddot{\mathbf{p}}_{k+2}^{\text{A}} - \pi_{\ddot{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_k))\|^2. \quad (44)$$

4.1 Incremental Attack Formulation

Due to feedback linearization, the manipulator's joint-space dynamics reduce to double integrators, which entails an *integrator vulnerability*. This vulnerability has been analyzed for sensor *Bias Injection Attacks* (BIAs), i.e., constant sensor injections, in linear systems (Tosun et al., 2025). Although the joint-task mapping is nonlinear, and we consider the more general *False Data Injection Attacks* (FDIAs) we argue that effective FDIAs still exploit the integrator vulnerability *locally*. For this reason, we model the attack *incrementally* as:

$$\mathbf{a}_k = \begin{bmatrix} (\mathbf{a}_k)_q \\ (\mathbf{a}_k)_{\dot{q}} \end{bmatrix} = \begin{bmatrix} (\mathbf{a}_{k-1})_q + \Delta_k \\ \frac{\Delta_k}{T_s} \end{bmatrix} = \begin{bmatrix} (\mathbf{a}_{k-1})_q \\ \mathbf{0} \end{bmatrix} + \mathbf{M} \Delta_k, \quad (45)$$

where $(\cdot)_q$, $(\cdot)_{\dot{q}}$ extract respectively the position and velocity components of the attack vector, $\Delta_k \in \mathbb{R}^n$ is an increment, and $\mathbf{M} = [\mathbf{I}_n \ \frac{1}{T_s} \mathbf{I}_n]^\top \in \mathbb{R}^{2n \times n}$ maps the joint increment to the full sensor space. This parameterizes an FDIA as deviations from a baseline BIA and it provides a convenient model for gradient-based synthesis.

A first-order expansion of $\pi_{\bar{\mathbf{p}}}(\mathcal{Z}_{k,2}(\cdot))$ around \mathbf{a}_{k-1} , gives:

$$\dot{\bar{\mathbf{p}}}_{k+2}^{\text{SIM}} \approx \pi_{\bar{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_{k-1})) + \mathbf{G}_k \Delta_k, \quad (46)$$

with $\ddot{\bar{\mathbf{p}}}_{k+2}^{\text{SIM}}$ as in (43) and the Jacobian matrix

$$\mathbf{G}_k = \mathbf{Z}_k \mathbf{M} \in \mathbb{R}^{3 \times n}, \quad \text{with} \quad (47)$$

$$\mathbf{Z}_k = \left. \frac{\partial}{\partial \mathbf{a}} \pi_{\bar{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a})) \right|_{\mathbf{a}=\mathbf{a}_{k-1}} \in \mathbb{R}^{3 \times 2n} \quad (48)$$

is computed numerically via a central-difference scheme. This \mathbf{G}_k provides the local sensitivity needed to differentiate the objective function in (44).

Substituting the linear approximation from (46) into the objective (44), through the relation in (43), and introducing a regularization term, the quadratic cost function is obtained. Specifically, the attack increment Δ_k minimizes

$$\frac{1}{2} \|\mathbf{G}_k \Delta_k - \ddot{\bar{\mathbf{p}}}_{k+2}^{\text{A}} + \pi_{\bar{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_{k-1}))\|^2 + \frac{\zeta}{2} \|\Delta_k\|^2, \quad (49)$$

where $\zeta > 0$ penalizes large increments.

The adversary defines $\ddot{\bar{\mathbf{p}}}_{k+2}^{\text{A}}$ in (44) via a one-step-ahead PD law to counteract the drift predicted for $k+1$ under $\Delta_k=0$ (i.e., using \mathbf{a}_{k-1}) as:

$$\ddot{\bar{\mathbf{p}}}_{k+2}^{\text{A}} = \mathbf{K}_p^{\text{A}} (\bar{\mathbf{p}}_{k+1}^{\text{A}} - \mathbf{p}_{k+1}^{\text{SIM}}) + \mathbf{K}_d^{\text{A}} (\dot{\bar{\mathbf{p}}}_{k+1}^{\text{A}} - \dot{\mathbf{p}}_{k+1}^{\text{SIM}}), \quad (50)$$

where

$$\dot{\bar{\mathbf{p}}}_{k+1}^{\text{SIM}} = \pi_{\bar{\mathbf{p}}}(\mathcal{Z}_{k,1}(\mathbf{a}_{k-1})), \quad (51)$$

$$\ddot{\bar{\mathbf{p}}}_{k+1}^{\text{SIM}} = \pi_{\bar{\mathbf{p}}}(\mathcal{Z}_{k,1}(\mathbf{a}_{k-1})), \quad (52)$$

are the one-step-ahead predicted position and velocity, and \mathbf{K}_p^{A} , \mathbf{K}_d^{A} are constant PD gains.

4.2 χ^2 Stealth Constraint

Due to the incremental modeling of the attack, the ADS residual of (10) is modeled as

$$\mathbf{r}_k = \mathbf{M} \Delta_k + \mathbf{c}_k, \quad (53)$$

where \mathbf{M} is defined in Section 4.1, and $\mathbf{c}_k \in \mathbb{R}^{2n}$ denotes the baseline innovation, defined as

$$\mathbf{c}_k = (\mathbf{y}_k + \mathbf{a}_{k-1} - \hat{\mathbf{y}}_k). \quad (54)$$

The stealth constraint becomes

$$(\mathbf{M} \Delta_k + \mathbf{c}_k)^\top \Sigma^{-1} (\mathbf{M} \Delta_k + \mathbf{c}_k) \leq \tau. \quad (55)$$

4.3 QCQP Formulation

Considering the objective function in (49) and expanding the constraint of (55), the problem reduces to the QCQP

$$\begin{aligned} \min_{\Delta_k} & \frac{1}{2} \Delta_k^\top \mathbf{H} \Delta_k + \mathbf{g}^\top \Delta_k \\ \text{s.t.} & (\mathbf{M} \Delta_k)^\top \mathbf{O} (\mathbf{M} \Delta_k) + \mathbf{b}^\top (\mathbf{M} \Delta_k) + c \leq 0, \end{aligned} \quad (56)$$

Table 1. Scenarios definitions.

Scenario	Task		Defence		
	Nominal	Attacker	χ^2	Damp.	Man. red.
A1	Circle	-	□	□	□
A2	Circle	-	□	■	□
A3	Circle	-	□	■	■
B1	Still	Circle	□	□	□
B2	Still	Circle	■	□	□
B3	Still	Circle	■	■	□
B4	Still	Circle	■	■	■

with parameters

$$\mathbf{H} = \mathbf{G}_k^\top \mathbf{G}_k + \zeta \mathbf{I},$$

$$\mathbf{g} = -\mathbf{G}_k^\top \left(\ddot{\bar{\mathbf{p}}}_{k+2}^{\text{A}} - \pi_{\bar{\mathbf{p}}}(\mathcal{Z}_{k,2}(\mathbf{a}_{k-1})) \right),$$

$$\mathbf{O} = \Sigma^{-1},$$

$$\mathbf{b} = 2\Sigma^{-1} \mathbf{c}_k,$$

$$c = \mathbf{c}_k^\top \Sigma^{-1} \mathbf{c}_k - \tau.$$

Theorem 5. (Convexity of adversary's QCQP). Problem (56) is convex and admits a unique global minimizer Δ_k^* .

Proof. The cost Hessian $\mathbf{H} = \mathbf{G}_k^\top \mathbf{G}_k + \zeta \mathbf{I} \succ \zeta \mathbf{I} \succ 0$, so the objective is strictly convex. The constraint Hessian $\mathbf{M}^\top \mathbf{O} \mathbf{M} = \mathbf{M}^\top \Sigma^{-1} \mathbf{M} \succ 0$, yielding a convex (ellipsoidal) feasible set. A strictly convex objective over a convex feasible set guarantees existence and uniqueness of the global minimizer Δ_k^* . \square

The attack evolves incrementally as in (45)

$$\mathbf{a}_k = \mathbf{a}_{k-1} + \Delta_k^*, \quad (57)$$

initialized with $\mathbf{a}_0 = \Delta_0^*$. At each step, the adversary run internal simulations to compute \mathbf{G}_k , $\mathbf{p}_{k+1}^{\text{SIM}}$ and $\dot{\mathbf{p}}_{k+1}^{\text{SIM}}$, then solves (56) to determine Δ_k^* , $k \in [0, T-1]$.

5. SIMULATION RESULTS

5.1 Experimental Setup

We define the following tasks:

Still Task: maintain the fixed hand position $\bar{\mathbf{p}}^{\text{task}} = [0.65, 0.17, 0.4]^\top$ m and RPY orientation $\theta^{\text{task}} = [3.05, -0.68, 1.26]$ rad.

Circle Task: Starting from the pose of the Still Task, perform a semi-circled trajectory with center in $\mathbf{0}$ (robot base) and final point $[-0.23, 0.26, 0.45]^\top$ m (quintic polynomials with zero boundary velocity/acceleration). Orientation is not controlled.

Table 1 defines Scenarios, indicating the tasks (nominal and for the attacker) and what defense system(s) are in place. Scenarios A1–A3 and B1–B4 represents respectively the absence (\mathcal{H}_0) and presence of FDIA as in (57).

We introduce the following metrics. Let $\bar{\mathbf{p}}_k$ and $\mathbf{p}_k \in \mathbb{R}^3$ denote the nominal and actual hand trajectories. Then:

$$\text{RMS}_{k=0}^{T-1}(\bar{\mathbf{p}}_k, \mathbf{p}_k) := \left(\frac{1}{T} \sum_{k=0}^{T-1} \|\bar{\mathbf{p}}_k - \mathbf{p}_k\|_2^2 \right)^{1/2},$$

$$\text{MND}_{k=0}^{T-1}(\bar{\mathbf{p}}_k, \mathbf{p}_k) := \max_{0 \leq k \leq T-1} \|\bar{\mathbf{p}}_k - \mathbf{p}_k\|_2.$$

The total absolute power of the virtual masses system is:

$$|P_k^{\text{tot}}| = \sum_{j=1}^n |u_{k,j} \cdot \dot{q}_{k,j}|.$$

The virtual kinetic energy respectively injected and dissipated by the defence damping of (26) is:

$$\mathcal{E}_{\text{inj}} = \sum_{k=0}^{T-1} \max(0, \dot{\mathbf{q}}_k^\top \mathbf{u}_k^{\text{d}}) T_s, \quad \mathcal{E}_{\text{diss}} = \sum_{k=0}^{T-1} |\min(0, \dot{\mathbf{q}}_k^\top \mathbf{u}_k^{\text{d}})| T_s.$$

The plant is a simulated 7-DOF Kinova Gen3 manipulator without joints limit, discretized with sampling time $T_s = 3 \cdot 10^{-3}$ s, and acceleration limits $\mathbf{u}_{\text{max}} = -\mathbf{u}_{\text{min}} = [1, 1, 1, 1, 10, 10, 10]$. Both joint positions and velocities are sensed, resulting in a measurement vector $\mathbf{y}_k \in \mathbb{R}^p$ with $p = 14$, sensing covariance $\mathbf{R} = 2.5 \cdot 10^{-9} \mathbf{I}_p$, and process noise covariance $\mathbf{Q} = \mathbf{Q}_{\text{base}} \otimes (q_c \mathbf{I}_n)$, where $\mathbf{Q}_{\text{base}} = \begin{bmatrix} T_s^3/3 & T_s^2/2 \\ T_s^2/2 & T_s \end{bmatrix}$ is the exact discretization of the Continuous White Noise Acceleration (CWNA) model $q_c = 10^{-5} \text{ rad}^2/\text{s}^3$ is the process noise intensity, and \otimes is the Kronecker product.

The χ^2 detector of Section 2.3 is calibrated for an ARL of 10^9 samples (833 hours), realized by threshold $\tau = 71.57$ for (11). The proposed adaptive damping defense of Section 3.2 is tuned with: $F = 3$, $\gamma = 6$, $\beta = 0.01$, $\psi = (1 - 0.01)$, determining $z_x = 29.14$. For Theorem 3, we have formal guarantees that, on average and under \mathcal{H}_0 , a single sample of the control command deviates from 1% of the nominal power every 100 samples. For the proposed manipulability reduction defense of Section 3.3 we set: $\alpha = 5 \cdot 10^3$, $\mathbf{D} = \text{diag}(10, 10, 10, 10, 1, 1, 1)$, and $\nu_{\text{max}} = 50$. All tasks has a duration of $T = 5000$ samples (15 s).

5.2 Comparative Analysis and Discussion

The proposed defense adds negligible computational overhead. The QCQP for the attacker (external to the defense) is solved in approximately 40 ms/step using a standard solver (Gurobi). The key performance metrics for all evaluated scenarios are summarized in Table 2. For scenarios A1–A3, the reported values are averaged over 100 Monte Carlo simulations. Figure 1 illustrates the time-series signals for a representative single run of each scenario.

Scenario A1 provides the reference nominal performance for the execution of the circle task under \mathcal{H}_0 . The proposed active defenses (A2, A3) introduce a negligible disturbance to the primary task, with hand displacement inaccuracy increasing by only $\approx 10^{-5}$ m. Compared to A1, the null-space motion required in A3 increases the mean and maximum power consumption ($|P_k^{\text{tot}}|$) by +950.3% and +2089.3%, respectively.

Scenario B1 provides the reference performance for the circle task compromised by FDIA to the sensors. The attack can steer the end-effector with relatively high accuracy. The presence of the χ^2 detector alone (B2) proves ineffective, as z_k remains well below the threshold τ . With the introduction of anomaly-aware virtual damping (B3, B4), the attacker faces a fundamental trade-off: increasing the injection raises the dissipation rate, while reducing the injection makes the stealthy attack less effective. Simultaneously, increasing the injection is bounded by the activa-

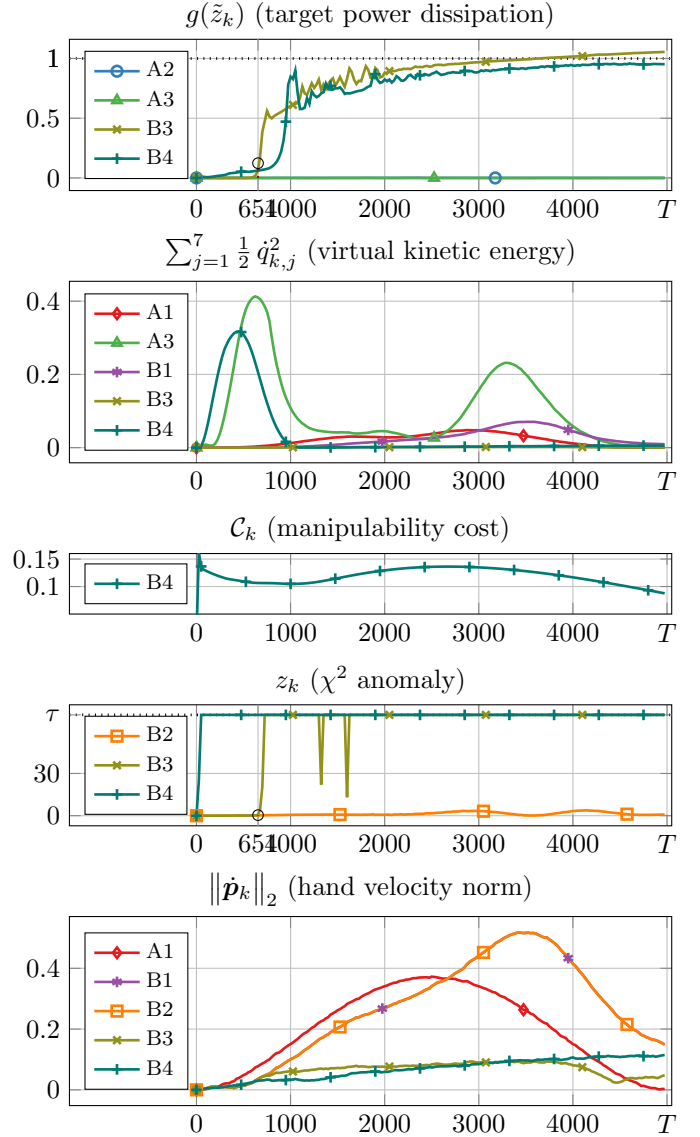


Fig. 1. Signals in the Scenarios as the sampling time k varies in $[0, T - 1]$.

tion of the stealthiness constraint in (55) enforcing $z_k < \tau$. Overall, the defense significantly reduces the kinetic energy of the double mass integrator, and the attack causes $g(\tilde{z}_k)$ to progressively approach 1, representing the limit above which the system becomes *strictly passive*. With the introduction of manipulability reduction (B4), a larger anomaly at the joint level is required to produce the same desired hand motion; this anticipates the activation of the stealthiness constraint and further decreases the impact of the attack. Notably, all attack impact metrics are reduced, e.g., $\text{RMS}_{k=0}^{T-1}(\tilde{\mathbf{p}}_k, \mathbf{p}_k)$ by -28.1% , $\text{RMS}_{k=0}^{T-1}(\tilde{\mathbf{p}}_k^{\text{A}}, \mathbf{p}_k)$ by $+14.2\%$, and the hand's path length by -14.8% . However, the mean and max of $|P_k^{\text{tot}}|$ increase by $+5801.6\%$ and $+1967.7\%$, respectively. As \tilde{z}_k increases, the virtual damping active defense (B3, B4) drives the actuators for joints q_2, q_3, q_4 toward their saturation limits, since the induced virtual friction progressively fills the actuator headroom reserved for defense in (25a) and (25b).

Across all experiments, the state deviation $\|\tilde{\mathbf{x}}_k - \mathbf{x}_k\|_2$ has an average of 0.11 and a maximum of 0.27. It is worth

Table 2. Comparative metrics results.

ID	Metric	A1	A2	A3	B1	B2	B3	B4
M1.1	$\text{RMS}_{k=0}^{T-1}(\hat{\mathbf{p}}_k, \mathbf{p}_k)$, [m]	$1.35 \cdot 10^{-4}$	$1.34 \cdot 10^{-4}$	$1.9 \cdot 10^{-4}$	0.58	0.58	0.27	0.19
M1.2	$\text{MND}_{k=0}^{T-1}(\hat{\mathbf{p}}_k, \mathbf{p}_k)$, [m]	$3.19 \cdot 10^{-4}$	$3.19 \cdot 10^{-4}$	$6.06 \cdot 10^{-4}$	0.88	0.88	0.45	0.39
M2.1	$\text{RMS}_{k=0}^{T-1}(\hat{\mathbf{p}}_k^A, \mathbf{p}_k)$, [m]	\emptyset	\emptyset	\emptyset	0.04	0.04	0.44	0.5
M2.2	$\text{MND}_{k=0}^{T-1}(\hat{\mathbf{p}}_k^A, \mathbf{p}_k)$, [m]	\emptyset	\emptyset	\emptyset	0.06	0.06	0.69	0.72
M3.1	mean $\ \dot{\mathbf{q}}_k\ _2$	0.17	0.17	0.32	0.19	0.19	0.04	0.16
M3.2	max $\ \dot{\mathbf{q}}_k\ _2$	0.31	0.31	0.82	0.38	0.38	0.07	0.8
M3.3	mean $ P_k^{\text{tot}} $	$9.73 \cdot 10^{-3}$	$9.73 \cdot 10^{-3}$	0.08	0.01	0.01	$8.05 \cdot 10^{-4}$	0.05
M3.4	max $ P_k^{\text{tot}} $	0.03	0.03	0.58	0.03	0.03	0.02	0.44
M3.5	path length of \mathbf{p}_k , [m]	1.22	1.22	1.22	1.24	1.24	0.48	0.41
M3.6	defence injected energy \mathcal{E}_{inj} , [J]	\emptyset	$2.9 \cdot 10^{-7}$	$4.84 \cdot 10^{-7}$	\emptyset	\emptyset	0.05	$3.88 \cdot 10^{-4}$
M3.7	defence dissipated energy $\mathcal{E}_{\text{diss}}$, [J]	\emptyset	$7.29 \cdot 10^{-5}$	$5.77 \cdot 10^{-4}$	\emptyset	\emptyset	0.14	0.11

discussing the impact of this deviation under \mathcal{H}_0 (A1, A2, A3), i.e., when $\hat{\mathbf{x}}_k \approx \mathbf{x}_k$. For the manipulability reduction defense, such deviation does not interfere with the primary task, because in (38) we project onto the null-space of $\hat{\mathbf{J}}_k \approx \mathbf{J}_k$. Regarding the damping defense, an excessive deviation could invalidate Assumption 5 and break the stability established in Theorem 3, i.e., if $\mathcal{E}_{\text{inj}} > \mathcal{E}_{\text{diss}}$. In our experiments, $\mathcal{E}_{\text{inj}} \ll \mathcal{E}_{\text{diss}}$; hence, stability is far from being compromised. We acknowledge that this phenomenon establishes an upper bound for the re-synchronization duration interval in (12).

6. CONCLUSIONS

This paper addressed the resilience of robotic manipulators against finite-horizon False Data Injection Attacks (FDIAs) targeting sensors. Feedback linearization induces an *integrator vulnerability* that allows sensor corruption to remain undetected by χ^2 anomaly detectors while driving the end-effector off-task. To counter this threat, we introduced *anomaly-aware virtual damping and manipulability reduction* in the attack direction, that disincentivate anomaly injection as a function of an anomaly score derived from a measurement-free, actuation-projected predictor. We established two key guarantees: (i) probabilistic bounds on power 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, providing a principled benchmark against which to evaluate defenses. Simulation results on a 7-DOF manipulator confirmed that the proposed defenses substantially improves the resilience against FDIA to robotic arms compared to using a χ^2 detection alone.

REFERENCES

- De Luca, A. and Oriolo, G. (1991). Issues in acceleration resolution of robot redundancy. In *IFAC Symp. Rob. Contr. (SYROCO)*, 93–98.
- Ding, D. et al. (2018). A survey on security control and attack detection for industrial cyber-physical systems. *Neurocomputing*, 275, 1674–1683.
- Fawzi, H., Tabuada, P., and Diggavi, S. (2014). Secure estimation and control for cyber-physical systems under adversarial attacks. *IEEE Trans. Autom. Contr.*, 59(6).
- Gualandi, G. and Papadopoulos, A.V. (2026). From passive monitoring to active defense: Resilient control of manipulators under cyberattacks. In *IEEE International Conference on Robotics and Automation (ICRA)*.
- Guo, Z., Shi, D., Johansson, K.H., and Shi, L. (2017). Optimal linear cyber-attack on remote state estimation. *IEEE Trans. Control Netw. Syst.*, 4(1), 4–13.
- Humayed, A., Lin, J., Li, F., and Luo, B. (2017). Cyber-physical systems security—A Survey. *IEEE Internet Things J.*, 4(6), 1802–1831.
- Intriago, A. et al. (2024). Residual-based detection of attacks in cyber-physical inverter-based microgrids. *IEEE Trans. Power Syst.*, 39(2), 4020–4038.
- Liao, L.Z. and Shoemaker, C. (1991). Convergence in unconstrained discrete-time differential dynamic programming. *IEEE Trans. Autom. Contr.*, 36(6), 692–706.
- Mo, Y. and Sinopoli, B. (2009). Secure control against replay attacks. In *Allerton Conf. Comm., Contr. & Comp.*, 911–918.
- Murguia, C. and Ruths, J. (2016). CUSUM and chi-squared attack detection of compromised sensors. In *IEEE Conf. Control Appl. (CCA)*, 474–480.
- Sandberg, H., Gupta, V., and Johansson, K.H. (2022). Secure networked control systems. *Annual Review of Contr., Rob., & Auton. Syst.*, 5(1), 445–464.
- Siciliano, B. and Slotine, J.J. (1991). A general algorithm for managing multiple tasks in highly redundant robotic systems. In *Int. Conf. Adv. Rob. (ICAR)*, 1211–1216.
- Siciliano, B., Sciavicco, L., Villani, L., and Oriolo, G. (2009). *Robotics: Modelling, Planning and Control*. Advanced Textbooks in Control and Signal Processing. Springer London.
- Tosun, F.E., Teixeira, A.M.H., Dong, J., Ahlén, A., and Dey, S. (2025). Kullback-Leibler divergence-based observer design against sensor bias injection attacks in single-output systems. *IEEE Trans. Information Forensics and Security*, 20, 2763–2777.
- Tunga, R., Murguia, C., and Ruths, J. (2018). Tuning windowed chi-squared detectors for sensor attacks. In *Am. Contr. Conf. (ACC)*, 1752–1757.
- Ueda, J. and Blevins, J. (2024). Affine transformation-based perfectly undetectable false data injection attacks on remote manipulator kinematic control with attack detector. *IEEE Robot. Autom. Lett.*, 9(10), 8690–8697.
- Van der Schaft, A. (2000). *L2-gain and passivity techniques in nonlinear control*. Springer.