

# Autonomous Cyber Resilience via a Co-Evolutionary Arms Race within a Fortified Digital Twin Sandbox

Malikussaid

*School of Computing, Telkom University*  
Bandung, Indonesia  
malikussaid@student.telkomuniversity.ac.id

Sutiyo

*School of Computing, Telkom University*  
Bandung, Indonesia  
tioatmadja@telkomuniversity.ac.id

**Abstract**—The convergence of Information Technology (IT) and Operational Technology (OT) has created hyper-connected Industrial Control Systems (ICS), exposing critical infrastructure to a new class of adaptive, intelligent adversaries that render static defenses obsolete. Existing security paradigms often fail to address a foundational “Trinity of Trust,” comprising the fidelity of the system model, the integrity of synchronizing data, and the resilience of the analytical engine against sophisticated evasion. This paper introduces the Adversarial Resilience Co-evolution (ARC) framework, a method for achieving analytical resilience through an autonomous, closed-loop hardening process. ARC establishes a perpetual co-evolutionary arms race within the high-fidelity sandbox of a Fortified Secure Digital Twin (F-SCDT). A Deep Reinforcement Learning (DRL) agent, the “Red Agent,” is formalized and incentivized to autonomously discover stealthy, physically-plausible attack paths that maximize process disruption while evading detection. Concurrently, an ensemble-based “Blue Agent” defender is continuously hardened via adversarial training against the evolving threats discovered by its adversary. This co-evolutionary dynamic forces both agents to become progressively more sophisticated, enabling the system to autonomously probe and patch its own vulnerabilities. Experimental validation on both the Tennessee Eastman Process (TEP) and the Secure Water Treatment (SWaT) testbeds demonstrates the framework’s superior performance. A comprehensive ablation study, supported by extensive visualizations including Receiver Operating Characteristic (ROC) curves and SHapley Additive exPlanations (SHAP) plots, reveals that the co-evolutionary process itself is responsible for a significant performance increase in detecting novel attacks. By integrating Explainable Artificial Intelligence (XAI) to ensure operator trust and proposing a scalable Federated ARC (F-ARC) architecture, this work presents ARC not merely as an improvement, but as a necessary paradigm shift toward dynamic, self-improving security for the future of critical infrastructure.

**Keywords**—Digital twin, industrial control systems security, co-evolutionary algorithms, adversarial machine learning, anomaly detection

## I. THE IMPERATIVE FOR DYNAMIC RESILIENCE IN CRITICAL INFRASTRUCTURE

The operational landscape of global critical infrastructure has undergone a profound and irreversible transformation, built upon a brittle foundation. The fourth industrial revolution, or Industry 4.0 [1], represents a paradigm shift in industrial automation, driven by the mass integration of cyber-physical systems, the Industrial Internet of Things (IIoT) [2], and cloud computing. This convergence of Information Technology (IT) and Operational Technology (OT) has irrevocably dismantled the “air gap” that once served as the primary, albeit fragile, defense for the world’s critical infrastructure. The resulting hyper-connected ecosystem of Industrial Control Systems (ICS) [3] has unlocked

unprecedented efficiencies but has also created a catastrophically expanded and dangerously porous attack surface. This vulnerability is not a recent phenomenon or an accidental oversight but the direct and predictable culmination of decades of design choices that, justifiably at the time, prioritized operational reliability, predictability, and physical safety over security in what was assumed to be a trusted, physically isolated, and non-adversarial environment.

This fragility stems not only from outdated protocols but also from a deeper, philosophical conflict between two disparate operational paradigms. The IT sector operates on a model of rapid development and continuous patching, while the OT sector is governed by a culture of extreme risk aversion, where maintaining system availability and stability is the paramount objective. This fundamental mismatch in the philosophy of change management and risk assessment created a critical seam—a vulnerability not in code but in process and culture—which adversaries quickly learned to exploit. The introduction of Artificial Intelligence (AI) into this precarious ecosystem has acted as a powerful catalyst, igniting a dynamic and perpetual co-evolutionary arms race. This necessitates a fundamental shift in security thinking, moving away from static, prevention-focused models toward resilient, adaptive paradigms.

The lessons from seminal attacks reveal that a robust security framework for modern ICS must be built upon a central, motivating thesis. This paper proposes this thesis as the “Trinity of Trust,” a set of three interconnected and non-negotiable foundational requirements. A failure in any one of these pillars renders the entire security posture invalid and creates an exploitable vulnerability. It is the failure of previous approaches to address all three pillars holistically that necessitates a new paradigm.

### A. The Foundational Flaws of the OT Environment

For much of the 20th century, the design philosophy of ICS was rooted entirely in the principles of safety and reliability engineering. The primary concerns were preventing physical accidents, ensuring process uptime, and guaranteeing deterministic behavior, with security being an entirely alien concept. This historical context is critical; the insecurity of modern OT is not an accident but a core feature of its design legacy.

The bedrock of this insecurity lies in the communication protocols themselves. Legacy protocols such as Modbus [4], developed in 1979 for serial communication between Programmable Logic Controllers (PLCs), and Distributed Network Protocol 3 (DNP3) [5], dating from 1993, still form the digital backbone of countless industrial facilities today, from power grids to water treatment plants and manufacturing floors. These protocols were designed for efficiency and simplicity in trusted, serial-line environments. As such, they

were created without even the most basic security features that are now considered fundamental in IT networking:

- The protocols lack any authentication mechanism to verify the identity of the device or user sending a command. Any device on the network can send a command to any other device, and it will be accepted as legitimate.
- All communication is transmitted in cleartext without encryption, meaning an attacker with network access can not only read all data but also easily intercept and modify commands in transit without being detected.
- Beyond basic error checking (like a CRC), there are no cryptographic integrity checks to ensure that a command has not been altered.

The practical implication of these design choices is a “trust-by-default” environment. An attacker with access to the OT network can trivially craft a valid Modbus TCP packet to, for example, switch off a critical cooling pump, and the target PLC will obey without question. There is no concept of a privileged user or a trusted source; all commands are treated with equal legitimacy.

The first documented cyber-physical attack, the 2000 Maroochy Shire sewage spill in Queensland, Australia [6], was a stark warning of the consequences of this inbuilt trust. A disgruntled former employee, using a laptop, radio transmitter, and knowledge of the system’s vulnerabilities, was able to remotely take control of the sewage system, ultimately releasing one million liters of raw sewage into local parks and rivers. While this incident clearly demonstrated that the threat was real, it was widely treated as an outlier—an act of insider revenge, not a harbinger of a new class of external, sophisticated threats. This misinterpretation led to a prolonged period of complacency, allowing the fundamental insecurities of OT to become even more deeply entrenched as systems were increasingly networked.

Yet, the incident was widely framed as an insider threat problem rather than a systemic technological failure. This framing proved to be a critical missed opportunity for the industry to confront the “insecure-by-design” problem at its core, leading to a prolonged period of complacency where the fundamental insecurities of OT became even more deeply entrenched. This has created a “technical debt” of insecurity that has now come due, compelling organizations to protect systems that were never designed to be defended.

#### *B. Threat Evolution from Nuisance to National Security Imperative*

The threat landscape has evolved dramatically from isolated incidents to persistent, sophisticated campaigns orchestrated by a complex and increasingly interconnected ecosystem of actors. Reports from leading government and industry bodies consistently highlight the escalating frequency and sophistication of attacks targeting these environments. The Dragos 2023 Year in Review reported a 49.5% increase in ransomware attacks impacting industrial organizations, with many incidents leading to operational shutdowns. By 2024, this trend had accelerated, with Dragos reporting an 87% year-over-year increase in ransomware attacks on the industrial sector, with 80 distinct ransomware groups active. This reality underscores a grim truth: the defenders of critical infrastructure are no longer facing isolated hackers but

adaptive, well-resourced adversaries conducting persistent campaigns.

This evolution is intrinsically linked to global geopolitical tensions, which now directly fuel OT-centric cyber operations. State-sponsored actors from Russia, China, and Iran have been identified as major threats, each with distinct motivations and tactics.

- Volt Typhoon (also tracked as VOLTZITE) [7], a People’s Republic of China (PRC)-affiliated actor, has been a significant concern due to its focus on pre-positioning within U.S. critical infrastructure. This group specializes in “living-off-the-land” (LotL) techniques, using legitimate, built-in system tools to evade detection and maintain long-term persistence. Their goal appears to be establishing footholds for potential future disruptive or destructive attacks in the event of a major conflict, with observed activities including the exfiltration of sensitive OT network diagrams and operational procedures from compromised utilities.
- APT44 (also known as Sandworm or ELECTRUM), a unit of the Russian GRU, has demonstrated a tactical evolution in the context of the war in Ukraine. Initially known for disruptive attacks, the group has increasingly focused on espionage to support conventional military operations. It has employed destructive wiper malware like ACIDPOUR, which can target OT devices, and has collaborated with hacktivist personas such as CyberArmyofRussia\_Reborn to create a layer of plausible deniability for its operations.

This convergence of state actors, cybercriminals, and hacktivists has lowered the barrier to entry for impactful attacks. Less sophisticated groups can now leverage leaked tools or basic techniques to cause tangible disruptions. For example, the Fuxnet malware, revealed in 2024, demonstrated how even a rudimentary tool could disrupt industrial sensors. Simultaneously, hacktivist groups like CyberArmyofRussia\_Reborn have successfully targeted internet-exposed Human-Machine Interfaces (HMIs), proving that high technical sophistication is not a prerequisite for causing operational and psychological impact.

#### *C. Anatomy of Seminal ICS Attacks and Their Human Impact*

A deeper analysis of seminal ICS attacks reveals not only the technical depth of the vulnerabilities exploited but also a clear progression from simple nuisance to sophisticated, physically-destructive campaigns that increasingly target the human operator.

- Stuxnet (2010): The discovery of Stuxnet [8] marked a turning point. Far more than just malware, Stuxnet was a precision weapon targeting Siemens Step7 PLCs controlling uranium enrichment centrifuges in Iran. Its genius lay in its multi-stage attack chain. It exploited four separate zero-day vulnerabilities for propagation and privilege escalation. Once on the target network, it did not cause immediate failure. Instead, it subtly manipulated the rotational frequency of the centrifuges, inducing extreme mechanical stress and causing physical damage over time [9]. Crucially, it did so while replaying normal operational data to the

operators, creating a phantom reality where the HMI showed a stable process while the machinery was tearing itself apart. This attack represents a masterclass in exploiting the gap between the digital representation and the physical reality, effectively weaponizing the operator's trust in their own control system.

- **Ukrainian Power Grid (2015):** This attack [10] demonstrated a coordinated, multi-pronged assault. The attackers used spear-phishing emails to deploy the BlackEnergy 3 malware, gaining a foothold in the IT networks of three energy distribution companies. From there, they pivoted to the Supervisory Control and Data Acquisition (SCADA) network, systematically mapping the environment for months. The final attack was executed in minutes: operators were locked out of their own systems as they watched the attackers' cursors move on their screens, remotely opening breakers at multiple substations and causing widespread power outages for over 230,000 consumers. The attack also included components to flood telecommunication systems and wipe firmware to complicate recovery efforts. The psychological impact on the operators, forced to be helpless spectators to the destruction of their grid, cannot be overstated and highlights the human-centric nature of modern attacks.
- **Cyber Av3ngers Campaign (2023):** More recent campaigns show an evolution towards psychological warfare and systemic disruption. The 2023 campaign by the "Cyber Av3ngers," an Iranian Revolutionary Guard Corps (IRGC)-affiliated persona [11], against U.S. water and wastewater systems serves as a visceral case study. The attack, while not technically complex, was devastatingly effective because it targeted trust. The attackers exploited basic security hygiene failures, targeting Israeli-made Unitronics Vision Series PLCs that were insecurely exposed to the internet, often using default passwords. Upon gaining access, they defaced the HMI with hostile messages. For an operator, seeing their trusted interface replaced with "You have been hacked" transforms a technical problem into a moment of helplessness and terror, instantly eroding confidence in the system's integrity. This underscores a critical reality: modern adversaries are not just exploiting code; they are exploiting the human operator's cognitive and emotional state [12].

The progression of these landmark attacks reveals a sophisticated learning curve among adversaries. The focus has expanded from purely technical exploitation of a physical process (Stuxnet) to the disruption of the human response loop (Ukraine power grid), the subversion of defensive tools' trust in legitimate processes (Volt Typhoon), and finally, to direct psychological attacks on human operators and the public (Cyber Av3ngers). Adversaries have learned that compromising the human element—by eroding trust, inducing stress, and creating debilitating uncertainty—can be as effective, if not more so, than simply breaking a physical component. This evolution demands that defensive strategies move beyond purely technical controls to incorporate human factors [13] as a core tenet of a holistic security posture.

These attacks demonstrate that adversaries orchestrate complex campaigns that blend digital intrusion with a deep understanding of the target physical process and its human

operators. Static, signature-based defenses are utterly insufficient, as an adversary can be digitally "valid" by issuing commands within normal protocol specifications, yet be physically malicious.

#### *D. The Trinity of Trust as a Foundational Security Paradigm for Cyber-Physical Systems*

The lessons from these attacks reveal that a robust security framework for modern ICS must be built upon what this work terms a "Trinity of Trust," a set of three interconnected foundational requirements. A failure in any one of these pillars renders the entire security posture invalid.

1. **Model Fidelity:** The digital representation (e.g., digital twin) [14] used for analysis must be a high-fidelity replica of the physical asset. A model that fails to accurately capture the underlying physics and dynamics of the process is fundamentally untrustworthy. A critical vulnerability is "model drift," where the digital twin becomes less accurate over time due to physical factors like equipment wear and tear or process reconfiguration. An adversary can exploit this drift, executing an attack that the outdated model no longer recognizes as anomalous. The Stuxnet attack was a masterclass in exploiting such a fidelity gap; the operators' HMI model showed a stable process while the physical centrifuges were being destroyed. True fidelity requires a hybrid model that understands the physical laws governing the process [15] to detect non-intuitive, but physically anomalous, states. Maintaining this fidelity is an ongoing challenge, requiring continuous validation and recalibration to counter semantic drift and ensure the digital representation remains a reliable source of truth.
2. **Data Integrity:** The data stream connecting the physical world to the digital analysis engine must be verifiably authentic and tamper-proof. Even with a perfect model, the system is blind if its data feeds can be deceived. As demonstrated by Stuxnet and Cyber Av3ngers, an adversary who can manipulate this data pipeline can effectively "gaslight" the entire system. Securing this pipeline with cryptographic guarantees like Elliptic Curve Digital Signature Algorithm (ECDSA) signatures and immutable ledgers is a necessary, non-negotiable step. However, integrity must be considered holistically. A compromised endpoint sensor can truthfully sign and transmit malicious data, bypassing gateway-level security. To address this, a hardware root of trust at the sensor level is required. Physically Unclonable Functions (PUFs) [16] offer a powerful solution by generating unique, device-specific cryptographic keys from the inherent, random physical variations of the silicon itself. Because these keys are generated on-demand and never stored in non-volatile memory, they are virtually impossible to clone or extract, providing a much stronger guarantee of data origin than traditional methods.
3. **Analytical Resilience:** The analytical engine—the anomaly detection models themselves—must be resilient to an intelligent adversary who actively seeks to evade it. Having a perfect model and pristine data is useless if the brain of the security system is easily outsmarted. Standard machine learning models trained

on historical data are fundamentally brittle, perfectly prepared for yesterday’s war but easily bypassed by a novel attack [17]. For example, a “low-and-slow” attack, where an adversary gradually manipulates a sensor value over hours, can keep each individual data point within a “normal” region, evading a simple anomaly detector that lacks temporal context. True resilience requires a defense that actively anticipates and adapts to the adversary’s evolving tactics, hardening itself against threats it has not yet seen.

#### E. Research Contributions: The Adversarial Resilience Co-evolution (ARC) Framework

This paper argues that achieving Analytical Resilience in a dynamic threat environment is the central challenge for ICS security. While architectures can address fidelity and integrity, new methods are required to build a defense that learns and evolves. To this end, this work makes the following primary contributions:

- The paper proposes the *Adversarial Resilience Co-evolution (ARC) framework*, a process for the autonomous, closed-loop hardening of ICS defenses that transforms a digital twin from a passive monitor into an active, self-improving security system.
- The research *formalizes the problem* of discovering stealthy, physically-plausible ICS attacks within a Deep Reinforcement Learning (DRL) paradigm, featuring a novel reward function explicitly engineered to balance the competing objectives of physical disruption and adversarial evasion.
- The work conducts comprehensive *experimental validation* on multiple, distinct testbeds—the Tennessee Eastman Process (TEP) [18] and the Secure Water Treatment (SWaT) benchmark [19]—including a novel ablation study that quantitatively analyzes the source of performance gains.
- The paper introduces a technically-grounded *vision for scaling autonomous resilience* with Federated Learning (F-ARC) [20], providing a blueprint for an industry-wide, privacy-preserving defense ecosystem and analyzing its unique security challenges.

## II. STATE-OF-THE-ART IN ICS DEFENSE

To justify the novelty and necessity of the ARC framework, this section provides an expanded critical analysis of prior art, systematically deconstructing the limitations of existing approaches when viewed through the lens of the Trinity of Trust—specifically, the pillar of Analytical Resilience. The evolution of ICS defense can be broadly categorized into four paradigms, each building upon the last, yet each introducing its own set of limitations that necessitate the next leap forward.

#### A. The Failure of Static Defenses in a Dynamic World

The first generation of ICS security solutions were direct imports from the IT world, a pragmatic but flawed response to the initial wave of IT/OT convergence. This convergence was driven by clear business imperatives: the need for real-time production data for enterprise resource planning (ERP) systems, remote monitoring capabilities, and predictive maintenance analytics. In connecting previously isolated OT networks to corporate IT networks, organizations inadvertently shattered the “air gap” and exposed legacy

systems, which were never designed for an adversarial environment, to a world of new threats. The initial security response was to deploy what was known and available: perimeter firewalls with static, rule-based access control lists, and network Intrusion Detection Systems (IDS) that relied on signatures of known malware [21].

While these tools form a necessary, foundational layer of basic security hygiene, they are fundamentally reactive and architecturally mismatched for the realities of the modern OT threat landscape. Their core logic is predicated on identifying known malicious patterns or blocking unauthorized communication paths. However, they are incapable of understanding the *context* or *intent* behind legitimate commands. An attacker using a standard Modbus Write Single Coil command (Function Code 05) to disable a critical safety interlock, followed by a Write Single Register command (Function Code 06) to push a motor’s speed beyond its operational limits, will not trigger any malware signature. The traffic itself is syntactically correct and protocol-compliant, and will be passed by any basic firewall that allows Modbus communication. This approach fails the test of Analytical Resilience because it is predicated on a static, uncreative adversary and cannot comprehend malicious *intent* when it is expressed using legitimate *syntax*.

The most critical and increasingly prevalent blind spot of this paradigm is its inability to counter “living-off-the-land” (LotL) techniques [22]. Sophisticated adversaries, such as the state-sponsored group Volt Typhoon, have demonstrated a clear preference for deliberately avoiding the deployment of custom, signature-able malware. Instead, they abuse legitimate, built-in system tools and protocols—such as PowerShell for scripting, Windows Management Instrumentation (WMI) for system queries, netsh for network configuration, and valid industrial commands—to achieve their objectives. Because these actions leverage trusted and often whitelisted system components, they blend seamlessly with normal administrative and operational traffic. From the perspective of a signature-based IDS or a basic firewall, these activities are indistinguishable from the actions of a legitimate system administrator or engineer performing routine maintenance. This tactic renders the entire class of signature-based defense largely obsolete against modern, persistent threats, as there is simply no “bad” signature to detect.

#### B. Promise, Pitfalls, and Operational Naivete of Machine Learning for Anomaly Detection

Recognizing the stark limitations of static signatures, the research community and security industry turned to machine learning (ML) and data-driven anomaly detection [23]. This represented a conceptual leap forward, moving from identifying known “badness” to modeling a baseline of “goodness,” or normal behavior, and flagging any deviation from it.

Early efforts focused on “shallow” machine learning models like Support Vector Machines (SVMs) [24], Random Forests, and K-Nearest Neighbors [25]. While these models proved effective in some contexts, they suffer from two major drawbacks in the ICS domain. First, they often require extensive and brittle feature engineering, demanding that a domain expert manually select the most relevant sensor data and statistical features. Second, and more critically, they typically treat sensor readings as independent features at discrete points in time, failing to capture the crucial temporal

dependencies and long-term sequential patterns that define a complex industrial process. For example, a Random Forest might correctly classify a single pressure reading of ‘101 psi’ as normal, but it cannot recognize that a sequence of readings—‘101 psi’, then ‘102 psi’, then ‘103 psi’—over three hours constitutes a dangerous trend indicative of a slow leak. A model without temporal context is blind to such attacks.

The advent of deep learning brought more powerful tools, particularly recurrent neural networks like Long Short-Term Memory (LSTM) networks [26], Gated Recurrent Units (GRUs) [27], and Autoencoders [28]. These architectures are inherently designed to model temporal sequences and learn complex, non-linear correlations in high-dimensional sensor data without manual feature engineering, making them a much better fit for the ICS domain. However, a persistent pitfall in much of the academic literature is the evaluation of these models only against pre-defined, non-adversarial fault conditions or simplistic, simulated attacks. When trained solely on benign operational data, these models learn a compressed representation of “normalcy,” effectively mapping the system’s normal operational states to a high-dimensional surface known as a manifold. This makes them inherently naive and brittle when faced with an intelligent adversary who will not behave like a random fault. Such an adversary will actively probe the model to understand the shape of this manifold and then craft an attack vector that constitutes a path *along* this surface from a safe state to a dangerous one, never deviating far enough from the manifold at any single point to be flagged as an anomaly. This fools the detector into classifying a malicious state transition as benign.

Furthermore, the application of deep learning in a critical operational context introduces its own set of significant challenges:

- **The Black Box Problem:** The complex, deeply-nested, non-linear nature of deep neural networks makes their decisions notoriously difficult to interpret [29]. If a model flags an anomaly, it often cannot explain *why* it did so in a way a human operator can understand and trust. This lack of transparency is a major barrier to adoption. An operator faced with an un-explainable alert must choose between ignoring a potentially critical warning or executing a costly and disruptive emergency shutdown based on blind faith in an algorithm. This can lead to “alert fatigue” and a general erosion of trust, causing operators to disable or ignore the security system altogether, rendering it useless.
- **Data Quality and Concept Drift:** Real-world industrial data is often noisy, and models trained on it are susceptible to a phenomenon known as “concept drift” [30]. Industrial processes are not static; equipment degrades, recipes are updated, and environmental conditions change. A model trained on data from six months ago may no longer accurately represent the “normal” state of the plant today, leading to a decay in performance and a rise in false alarms. Without a mechanism for continuous adaptation, the model’s fidelity inevitably declines.
- **Data Poisoning Vulnerabilities:** An adversary with even limited access to the model’s training data can execute a data poisoning attack [31]. By injecting a small number of carefully crafted malicious samples into the training set, the adversary can degrade the

model’s overall performance or, more insidiously, create a backdoor that causes the model to systematically misclassify a specific attack vector that the adversary intends to use later.

These limitations underscore that while deep learning is a powerful tool, its naive application is not a panacea and can even introduce new, subtle risks into the security posture.

### C. The Mismatch of Adversarial Machine Learning in Cyber-Physical Contexts

The field of Adversarial Machine Learning (AML) [32] emerged specifically to address the naivete of standard ML models by explicitly training them to be robust against malicious inputs designed to cause misclassification. However, a critical methodological error, prevalent in early research, is to directly port techniques and threat models from the domain of computer vision (CV), where AML was pioneered, to the cyber-physical domain of ICS.

The canonical threat model in computer vision is based on an  $L_\infty$  norm [33], which assumes the attacker can add a small, uniform, and often human-imperceptible amount of noise to *all* features (i.e., every pixel in an image). This is a physically implausible threat model for an industrial control system. It is not feasible for an attacker to subtly manipulate every single sensor reading across an entire plant simultaneously in a perfectly coordinated fashion. An attacker who compromises a system is far more likely to gain control over a limited number of sensors or actuators.

A far more realistic threat model, which this work adopts, involves an attacker compromising a small, discrete number of sensors or actuators. This corresponds to an attack that minimizes the  $L_0$  norm—the number of perturbed features [34]. To put this in an intuitive analogy: an  $L_\infty$  attack is like lightly dusting an entire landscape with a fine layer of snow, subtly changing everything at once. An  $L_0$  attack is like targeting a single house with a well-aimed snowball. In ICS, attackers throw snowballs; they do not change the global weather. The Jacobian-based Saliency Map Attack (JSMA) [35] is a classic and powerful  $L_0$  attack that serves as an excellent proxy for this type of focused, physically plausible attacker. JSMA works by calculating the Jacobian matrix of the model’s output with respect to its inputs, which effectively creates a “saliency map” indicating which input features have the most influence on the output class. It then iteratively modifies the small number of features that have the most significant impact, allowing it to efficiently craft a misclassification with minimal, targeted changes. By incorporating a physically-grounded  $L_0$  attack model like JSMA into the ARC hardening process, the framework trains its defender against a more realistic and dangerous class of adversary than those considered by standard AML techniques ported from computer vision.

### D. The Limits of Game Theory and Open-Loop Co-evolution

The most advanced thinking in the field views ICS security as a dynamic, strategic game between an attacker and a defender. This has led to two main lines of research: formal game-theoretic models and co-evolutionary algorithms.

*Game-theoretic models* [36], particularly those based on Stackelberg security games [37], provide a powerful formal framework for analyzing strategic interactions and resource allocation. In these models, a defender (the “leader”) commits

to a defensive strategy, and an attacker (the “follower”) observes this strategy and chooses a best response. While intellectually elegant, these models often remain at a high level of abstraction. Their practical application to a complex, real-world ICS is hampered by the immense difficulty of accurately defining the action spaces and, most critically, the payoff matrices for both players. Quantifying the precise “cost” of a sensor compromise or the “value” of a successful but partial disruption is often intractable. For an attacker, is the payoff measured in dollars of lost production, the severity of physical damage, or the psychological impact on operators? These values are highly subjective, context-dependent, and difficult to quantify, which is a major reason game theory models often remain in the realm of theoretical analysis rather than becoming concrete, implementable algorithms for real-time defense.

Other research has proposed *iterative or co-evolutionary schemes* [38], which come conceptually closer to the goal of an adaptive defense. Yet, they often lack the core feedback loop that defines a true, tightly-coupled arms race. For instance, many frameworks describe an “open-loop” process where an attacker generates a static set of new attacks, and a defender is subsequently retrained on this set. However, in these schemes, the attacker’s generation process is not influenced by the defender’s *updated* state. This is akin to a boxer training against a static video of a previous opponent’s fight; they might get very good at countering that *one* specific

style, but they are not prepared for a live opponent who adapts in real-time. This one-way adaptation is not a truly reciprocal co-evolution.

The ARC framework, in contrast, explicitly formalizes and implements this dynamic, closed-loop coupling. It is analogous to a live sparring session. The Red Agent’s reward function is directly penalized by the *current* Blue Agent’s real-time anomaly score, explicitly rewarding it for fooling the updated defender. In turn, the Blue Agent is explicitly hardened against the *emergent* strategies that were successful in the previous epoch. This creates a perpetual, reciprocal arms race where both agents are forced to become progressively more sophisticated. The goal of this process is to foster “emergent behavior”—the discovery of attack vectors and defensive strategies that were not pre-programmed or anticipated by the human designers at the outset. It is this continuous, closed-loop process that drives the autonomous discovery and patching of complex system vulnerabilities that would be missed by any static or open-loop approach.

Table I provides a critical analysis of state-of-the-art ICS defense frameworks, systematically positioning the ARC framework against existing classes of solutions and highlighting its unique combination of autonomous, adaptive hardening and explainability.

TABLE I. CRITICAL ANALYSIS OF STATE-OF-THE-ART ICS DEFENSE FRAMEWORKS

Framework / Approach	Detects Zero-Days?	Resilient to Adaptive Adversary?	Autonomous Hardening?	Provides Explanations (XAI)?	Scalability Model	Key Limitation
ARC Framework (This Work)	Yes (via DRL)	Yes (via Co-Evolution)	Yes (Closed-Loop)	Yes (SHAP)	Federated (F-ARC)	High initial DT modeling cost
Signature-Based IDS	No	No	No	No	Monolithic	Reactive; cannot detect novel or “living-off-the-land” attacks.
Shallow ML (e.g., SVM)	Limited	No	No	No	Monolithic	Fails to model temporal dependencies; brittle.
Naive Deep Learning (LSTM/AE)	Limited	No	No	No	Monolithic	Vulnerable to adversarial examples that mimic normal data.
CV-based Adversarial Training	No	Limited (wrong threat model)	No	No	Monolithic	Uses physically implausible $L_\infty$ threat models.
Abstract Game Theory	Yes (in theory)	Yes (in theory)	No	No	Monolithic	Lacks concrete, implementable algorithms for real systems.
Blockchain for Data Provenance	No	N/A (Integrity tool)	No	No	Distributed	Addresses data integrity only; performance overhead.

### III. F-SCDT: A HIGH-FIDELITY SANDBOX FOR ADVERSARIAL SIMULATION

The ARC framework’s co-evolutionary arms race requires a specific environment: a high-fidelity, securely synchronized digital twin that can serve as a realistic and trustworthy “sparring gym” for the attacker and defender agents. A low-fidelity twin would lead to the discovery of irrelevant vulnerabilities that do not exist in the physical world, while an insecure twin would allow the agents to “cheat,” undermining the entire training process. This section details the architecture of the Fortified Secure Digital Twin (F-SCDT), which establishes the foundational pillars of Model Fidelity and Data Integrity, providing the necessary ground truth for the subsequent adversarial training.

#### A. The Imperative of Hybrid Physics-Informed and Data-Driven Models to Achieve Model Fidelity

To achieve the high fidelity required for realistic adversarial simulation, a purely physics-based or purely data-driven model is insufficient. Physics-based models, while providing a strong theoretical foundation, often fail to capture unmodeled dynamics, equipment degradation, and the subtle, unique behaviors of a specific physical plant. Conversely, purely data-driven models, while excellent at learning complex correlations from data, lack physical grounding, making them prone to un-physically-plausible predictions when faced with out-of-distribution inputs—precisely the kind of inputs an adversary seeks to create.

The F-SCDT therefore employs a hybrid model that fuses a first-principles physics model with a data-driven model that learns the residual error, combining the strengths of both approaches [39]. For the Tennessee Eastman Process (TEP) reactor, modeled as a non-isothermal Continuous Stirred-

Tank Reactor (CSTR), the governing ordinary differential equations (ODEs) provide a robust physical baseline.

Equation (1) presents the mass balance on reactant A:

$$\frac{dC_A}{dt} = \frac{F}{V}(C_{Af} - C_A) - k_0 e^{-\frac{E}{RT}} C_A \quad (1)$$

where  $C_A$  is the concentration of reactant A in the reactor (mol/L),  $F$  is the volumetric flow rate into the reactor (L/min),  $V$  is the reactor volume (L),  $C_{Af}$  is the feed concentration of reactant A (mol/L),  $k_0$  is the pre-exponential factor (1/min),  $E$  is the activation energy (J/mol),  $R$  is the universal gas constant (J/(mol·K)), and  $T$  is the reactor temperature (K).

Equation (2) shows the energy balance for the reactor:

$$\frac{dT}{dt} = \frac{F}{V}(T_f - T) + \frac{-\Delta H}{\rho C_p} k_0 e^{-\frac{E}{RT}} C_A - \frac{UA}{\rho C_p V}(T - T_c) \quad (2)$$

where  $T_f$  is the feed temperature (K),  $\Delta H$  is the heat of reaction (J/mol),  $\rho$  is the fluid density (kg/L),  $C_p$  is the specific heat capacity (J/(kg·K)),  $U$  is the overall heat transfer coefficient (J/(min·m<sup>2</sup>·K)),  $A$  is the heat transfer area (m<sup>2</sup>), and  $T_c$  is the coolant temperature (K). These equations form the physics-informed component,  $f(x, u, d)$ .

However, this model is never perfect. To capture the residual error,  $\epsilon$ , arising from unmodeled phenomena such as catalyst degradation, equipment wear, or minor process fluctuations, a data-driven model is trained to learn this non-linear error from historical operational data. The research chose a Gated Recurrent Unit (GRU) network over the more common Long Short-Term Memory (LSTM) network for this task. While both are effective at modeling temporal sequences, GRUs possess a simpler architecture with fewer parameters (two gates—a reset gate and an update gate—versus LSTM’s three gates—input, output, and forget). This results in comparable performance with lower computational overhead and faster training times, a critical advantage for the high-iteration training loops within the ARC framework where the digital twin must be simulated repeatedly.

The hybrid model is trained in two stages. First, the parameters of the physics model (like the heat transfer coefficient  $U$  or the activation energy  $E$ ) are estimated from steady-state historical data using optimization techniques such as non-linear least squares. This anchors the model in the plant’s general operating characteristics. Second, the GRU is trained on dynamic operational data to predict the residual error—the difference between the physics model’s output and the actual sensor readings. Equation (3) expresses the final, high-fidelity Hybrid Model:

$$\dot{x}_{\text{hybrid}} = f(x, u, d) + g_{\text{GRU}}(x, u, \theta) \quad (3)$$

where  $\dot{x}_{\text{hybrid}}$  represents the state derivatives from the hybrid model,  $f(x, u, d)$  is the physics-based component with states  $x$ , control inputs  $u$ , and disturbances  $d$ , while  $g_{\text{GRU}}(x, u, \theta)$  is the GRU-based correction with parameters  $\theta$ .

This hybrid approach ensures that the digital twin is not only grounded in physical laws but also adapts to the real-world imperfections of the system. This is crucial for

preventing “semantic drift,” where the twin’s representation of reality slowly diverges from the actual physical asset over time, and for ensuring the twin remains a trustworthy simulation environment for discovering realistic vulnerabilities.

### B. Ensuring Data Integrity Through a Verifiable Cyber-Physical Data Pipeline

Data Integrity is addressed by a verifiable data pipeline that provides strong guarantees of data authenticity and tamper-proofing from the sensor to the analysis engine. The threat model here considers data injection, modification, and replay attacks at both the sensor and network levels. The F-SCDT architecture combines a hardware data diode [40] for enforcing unidirectional data flow from the OT network to the IT network, and a permissioned blockchain [41] (e.g., Hyperledger Fabric) for creating an immutable audit trail.

An Industrial Internet of Things (IIoT) Gateway located on the OT network is responsible for batching sensor data. To ensure data origin authentication, this gateway signs the data batch using a private key securely stored in a Hardware Security Module (HSM) [42]. The cryptographic hash of the signed batch is then recorded as a transaction on the permissioned blockchain, creating a permanent and tamper-evident record. Finally, the signed data batch is transmitted through the data diode to the F-SCDT environment. A data diode is a hardware-based cybersecurity device that enforces one-way data flow using a physical separation (typically optical), making it physically impossible for data, malware, or commands to flow back from the less trusted IT network into the critical OT network.

A more advanced approach to ensuring integrity at the source involves the use of Physically Unclonable Functions (PUFs) [43] embedded within the sensors themselves. A PUF leverages minute, random variations in a chip’s physical microstructure, introduced during manufacturing, to generate a unique and unclonable digital fingerprint for that specific device. This fingerprint can be used to derive a cryptographic key that is never stored in non-volatile memory and is effectively part of the hardware’s intrinsic identity. A PUF-enabled sensor could thus sign its own data at the point of creation, providing a hardware root of trust that is resilient to cloning and physical tampering, offering a stronger guarantee of data origin than a gateway-level HSM.

Practical challenges remain, particularly the trade-off between the transaction throughput of the blockchain and the desired temporal granularity of the data. Small batch sizes increase cryptographic and network overhead, while large batch sizes reduce the forensic granularity available for incident analysis. Algorithm 1 outlines the conceptual data sealing process.

---

#### Algorithm 1: Data Provenance Sealing

---

**Require:**  $D_{\text{batch}}$ : A data batch from IIoT Gateway.

1: Procedure **SealDataBatch**( $D_{\text{batch}}$ )

#### Origin Authentication and Data Signing

2:  $k_{\text{priv}} \leftarrow \text{GetKeyFromHSM}()$  or  $\text{DeriveFromPUF}()$

3: **if**  $k_{\text{priv}}$  is unavailable **then**

4:     Log(“CRITICAL: Security module key unavailable”)

---

---

**Algorithm 1: Data Provenance Sealing**

---

```
5:   return FAILURE
6: end if
7:  $D_{signed} \leftarrow \text{SignECDSA}(D_{batch}, k_{priv})$ 
Immutable Forensic Record Creation
8:  $h_{batch} \leftarrow \text{SHA-256}(D_{signed})$ 
9: status, receipt  $\leftarrow \text{SubmitToLedger}(h_{batch})$ 
10: if status is CONNECTION_ERROR then
11:   Log("WARNING: Ledger connection failed,
      audit trail suspended")
12: else if receipt is INVALID then
13:   Log("WARNING: Blockchain transaction
      failed, audit trail incomplete")
14: end if
Secure Unidirectional Transfer
15: status  $\leftarrow \text{PublishToDiode}(D_{signed})$ 
16: if status is FAILURE then
17:   Log("CRITICAL: Data Diode publishing
      failed")
18:   return FAILURE
19: end if
20: return SUCCESS
```

---

Algorithm 1 presents the data provenance sealing process, which ensures the integrity and authenticity of sensor data batches before they are transmitted to the analysis environment. The algorithm first authenticates the data origin using digital signatures generated with keys from either a Hardware Security Module (HSM) or Physically Unclonable Functions (PUFs). It then creates an immutable record of the data by submitting a cryptographic hash to a permissioned blockchain ledger. Finally, it securely publishes the signed data through a hardware data diode that enforces one-way flow from the operational technology (OT) network to the information technology (IT) network.

#### IV. ARC FRAMEWORK AS A METHODOLOGY FOR CO-EVOLUTIONARY ARMS RACE

This section details the ARC framework, which formalizes the interaction between an attacker and a defender into a perpetual, automated arms race to achieve the third pillar of the Trinity of Trust: Analytical Resilience. This is accomplished by creating a closed-loop system where two agents, a "Red Agent" attacker and a "Blue Agent" defender, are forced to continuously adapt to one another.

##### A. The Red Agent: Autonomous Vulnerability Discovery via Deep Reinforcement Learning

The Red Agent is a Deep Reinforcement Learning (DRL) agent tasked with finding stealthy and physically plausible attack vectors within the F-SCDT. The problem is formalized as a Markov Decision Process (MDP) [44], which consists of a tuple  $(S, A, P, R, \gamma)$ , where  $S$  is the set of states,  $A$  is the set of actions,  $P$  is the state transition probability function,  $R$  is the reward function, and  $\gamma$  is the discount factor.

- The research selected Proximal Policy Optimization (PPO) [45] as the core DRL algorithm. While off-policy algorithms like Soft Actor-Critic (SAC) [46]

can be more sample-efficient, they often suffer from training instability, especially in environments with high stochasticity and noisy state transitions, which are characteristic of complex ICS simulations. PPO, an on-policy algorithm, offers a superior balance of sample efficiency, stability, and ease of implementation. Its defining feature is a clipped surrogate objective function, which constrains the size of policy updates in each training step. This prevents destructively large changes to the policy, which is critical for ensuring stable and reliable convergence when the agent is interacting with a high-fidelity, sensitive simulation environment like the F-SCDT.

- *State Space ( $S$ ):* The state  $s_t$  is a comprehensive vector containing all relevant process variables from the digital twin (e.g., temperatures, pressures, flow rates). Crucially, it also includes the anomaly score generated by the Blue Agent defender at the previous time step. This allows the Red Agent to perceive the defender's state and learn which of its actions are being detected, enabling it to adapt its strategy toward greater stealth.
- *Action Space ( $A$ ):* The action  $a_t$  is a continuous-valued vector representing manipulations to actuator setpoints (e.g., valve positions, pump speeds). To ensure physical plausibility, these actions are clipped to realistic operational ranges, preventing the agent from learning to execute physically impossible state jumps.
- The reward function is the core of the agent's intelligence, meticulously engineered to incentivize the discovery of attacks that achieve maximum physical disruption while remaining undetected. Equation (4) defines this reward function:

$$R(s_t, a_t) = w_{disrupt} \times Disruption(s_{t+1}) - w_{detect} \times DetectionScore(s_{t+1}) \quad (4)$$

where  $w_{disrupt}$  and  $w_{detect}$  are weighting factors that balance disruption and stealth objectives. The *Disruption* term is a weighted sum of normalized deviations from critical operational setpoints and safety limits, directly rewarding the agent for pushing the system toward an unsafe or inefficient state. The *DetectionScore* is the maximum anomaly score produced by any model in the Blue Agent's ensemble, providing a strong and direct penalty for being detected. The weights ( $w_{disrupt}, w_{detect}$ ) are not arbitrary; they are critical hyperparameters determined through a rigorous sensitivity analysis [47]. This process involves systematically varying the weights and observing the Red Agent's emergent behavior to find a balance that encourages sophisticated, multi-stage attacks over simple, brute-force disruption that would be easily detected.

##### B. The Blue Agent: An Adaptive Ensemble for a Moving Target Defense

The Blue Agent is not a single model but an ensemble of diverse anomaly detectors, creating a multi-faceted defensive surface that is inherently more difficult for an adversary to evade. This ensemble consists of:



1. An LSTM network [26], chosen for its proven ability to capture long-term temporal dependencies in sequential sensor data, making it effective against “low-and-slow” attacks.
2. An Autoencoder [28], which excels at learning a compressed representation of the system’s normal state and identifying deviations based on reconstruction error. This makes it sensitive to attacks that violate complex physical correlations between multiple sensors.
3. An Isolation Forest [48], a tree-based algorithm that is efficient and effective for general-purpose outlier detection, capable of identifying novel anomalies that may not fit the patterns learned by the neural network models.

This diversity embodies the principle of a moving target defense at the algorithmic level; an adversary must learn to simultaneously fool three different detection logics, a significantly harder task than defeating a single model. The hardening process uses the novel attack vectors discovered by the Red Agent, further diversified with  $L_0$ -norm perturbations from the JSMA attack, to create an augmented training set for continuously improving the ensemble.

### C. The Co-Evolutionary Loop

The ARC algorithm formalizes the perpetual, closed-loop arms race between the Red and Blue agents. A key challenge in this continual learning process is “catastrophic forgetting” [49], where the defensive model, in learning to counter new threats, forgets how to detect older ones. The framework mitigates this with two mechanisms within the Train\_Defender function: a replay buffer that stores all previously discovered attacks [50], and a balanced sampling strategy. Each augmented training batch is composed of 50% normal data, 20% newly discovered attacks from the DRL agent, 10% JSMA-diversified versions of these new attacks, and 20% randomly sampled attacks from the historical replay buffer. This ensures the defender trains on a comprehensive distribution of threats, maintaining its knowledge base while adapting to the latest emergent strategies. Furthermore, the inclusion of JSMA-generated attacks serves to diversify the training data beyond what the DRL agent might discover on its own, preventing the two agents from over-fitting to each other’s specific strategies and forcing the defender to generalize against a wider range of physically plausible  $L_0$ -norm attacks.

Algorithm 2 presents the Adversarial Resilience Co-evolution (ARC) process, which implements a closed-loop arms race between attacker and defender agents.

---

#### Algorithm 2: Adversarial Resilience Co-evolution (ARC)

---

##### Require:

1.  $D_0$ : Initial defender models
  2.  $Z_{normal}$ : Normal data
  3.  $Z_{fault}$ : Known fault data
  4.  $N_{epochs}$ : Total number of co-evolutionary epochs
  5.  $N_{A_{epochs}}$ : Number of training steps for attacker per epoch
  6.  $N_{D_{epochs}}$ : Number of training steps for defender per epoch
- 

---

#### Algorithm 2: Adversarial Resilience Co-evolution (ARC)

---

7.  $w_{disrupt}$ : Reward weight for disruption

8.  $w_{detect}$ : Reward weight for detection

##### Initialize:

1.  $A_0$ : Initial attacker DRL agent
2.  $Z_{attacks} \leftarrow \emptyset$ : Cumulative set of all generated attacks

1: **for** epoch = 1 to  $N_{epochs}$  **do**

##### Attacker (“Red Agent”) Training Phase

- 2:  $D_{current} \leftarrow D_{epoch-1}$
- 3: **function** Train\_Attacker( $A_{epoch-1}, D_{current}$ )
- 4:   **for** a\_step = 1 to  $N_{A_{epochs}}$  **do**
- 5:     Collect trajectory  $\tau$  in DT using policy  $\pi_{A_{epoch-1}}$
- 6:     **for each**  $(s_t, a_t)$  in  $\tau$
- 7:        $s_{t+1} \leftarrow DT(s_t, a_t)$
- 8:        $D_{score} \leftarrow D_{current}(s_{t+1})$
- 9:        $R_t \leftarrow w_{disrupt} \times Disruption(s_{t+1})$   
            $-w_{detect} \times D_{score}$
- 10:      Update policy  $\pi_{A_{epoch-1}}$  using PPO with rewards  $\{r_t\}$
- 11:    **end for**
- 12: **end for**
- 13:   **return** updated agent  $A_{epoch}$
- 14: **end function**
- 15:  $A_{epoch} \leftarrow \text{Train\_Attacker}(A_{epoch-1}, D_{current})$

##### Generate New Attack Dataset from the Trained Attacker

- 16:  $Z_{new} \leftarrow \text{Generate\_Attacks}(A_{epoch}, DT, \text{num\_samples})$
- 17:  $Z_{attacks} \leftarrow Z_{attacks} \cup Z_{new}$

##### Defender (“Blue Agent”) Hardening Phase

- 18: **function** Train\_Defender( $D_{epoch-1}, Z_{new}, Z_{attacks}$ )

*Diversify new attacks with  $L_0$  perturbations to improve generalization*

- 19:  $Z_{JSMA} \leftarrow Z_{JSMA}(Z_{new}, D_{epoch-1})$

*Create augmented training set with replay to prevent catastrophic forgetting*

- 20:  $Z_{old\_sample} \leftarrow \text{Sample}(\frac{Z_{attacks}}{Z_{new}}, \text{ratio} = 0.2)$
  - 21:  $Z_{aug} \leftarrow Z_{normal} \cup Z_{fault} \cup Z_{new} \cup Z_{JSMA} \cup Z_{old\_sample}$
  - 22:   **for** d\_step = 1 to  $N_{D_{epochs}}$  **do**
  - 23:     Sample mini-batch from  $Z_{aug}$
  - 24:     Update weights of  $D_{epoch-1}$  via backpropagation on ensemble loss
  - 25:   **end for**
-

**Algorithm 2: Adversarial Resilience Co-evolution (ARC)**

```

26:   return updated defender  $D_{epoch}$ 
27: end function
28:    $D_{epoch} \leftarrow \text{Train\_Defender}(D_{epoch-1}, Z_{new},$ 
     $Z_{attacks})$ 
29: end for
30: return  $D_{N_{epochs}}$ 

```

## V. EXPERIMENTAL VALIDATION AND ANALYSIS OF EMERGENT VULNERABILITIES

This section presents the experimental validation of the ARC framework. The goal is not merely to demonstrate superiority over a static baseline but to show ARC’s utility as an automated process for discovering, analyzing, and patching complex, emergent vulnerabilities that would otherwise remain hidden.

### A. Experimental Setup

Validation was conducted on two distinct and well-regarded ICS security testbeds: the Tennessee Eastman Process (TEP) [18] and the Secure Water Treatment (SWaT) testbed [51]. The TEP is a complex, nonlinear chemical process *simulation* that is a standard benchmark for process control and security research, allowing for safe experimentation with highly disruptive attacks. SWaT is an operational *physical testbed*—a scaled-down but fully functional water treatment plant—that provides validation on real hardware, networks, and physical dynamics. Using both a complex simulation and a real-world testbed demonstrates the generalizability and practical applicability of the ARC process.

The baseline defender model,  $D_0$ , was trained on a comprehensive dataset of normal operational data and a set of standard, pre-defined fault scenarios (e.g., single sensor failures, valve malfunctions) before the ARC co-evolutionary process began. This represents a typical, state-of-the-practice anomaly detection system.

### B. Main Results and Ablation Study

The ARC-hardened model demonstrates significantly superior performance, particularly in detecting the novel, stealthy attacks discovered by the Red Agent during the co-evolutionary process. As shown in Table II, the  $F_1$ -Score for detecting a DRL-Discovered Stealth Attack improved from 0.65 for the baseline model to 0.89 for the ARC-hardened model. Furthermore, the detection latency was reduced from over 1200 seconds (with many attacks being missed entirely) to 210 seconds. This reduction is operationally critical, as it can be the difference between isolating a compromised segment of the plant and suffering a full-scale, cascading failure. Table II presents key performance metrics of the ARC framework versus baseline on TEP.

TABLE II. KEY PERFORMANCE METRICS OF THE ARC FRAMEWORK VS. BASELINE ON TEP

Scenario	Model	$F_1$ -Score	Detection Latency (s)
Data Replay Attack	ARC-Hardened	0.93	65
	Baseline ( $D_0$ )	0.73	450
DRL-Discovered Stealth Attack	ARC-Hardened	0.89	210

Scenario	Model	$F_1$ -Score	Detection Latency (s)
	Baseline ( $D_0$ )	0.65	> 1200 (often missed)

To deconstruct the sources of this performance gain, a rigorous ablation study [52] was conducted, the results of which are presented in Table III. This study systematically removes components from the full ARC framework to quantify their individual contributions. The results show that while each component of the defensive ensemble contributes to overall performance, the single most critical factor is the ARC training process itself. Removing the co-evolutionary hardening process while keeping the full defensive ensemble results in a massive 27.0% degradation in  $F_1$ -Score on the DRL-discovered attack, dropping performance back to the baseline level. This provides strong quantitative evidence that it is the co-evolutionary dynamic—the arms race—that forges the model’s resilience to novel, intelligent threats. Table III shows the ablation study of detection engine components.

TABLE III. ABLATION STUDY OF DETECTION ENGINE COMPONENTS

Configuration	$F_1$ -Score (on DRL-Discovered Attack)	Performance Degradation
ARC Framework (Full System)	0.89	-
ARC Framework without LSTM	0.82	-7.9%
ARC Framework without Autoencoder	0.80	-10.1%
Full Ensemble without ARC Training	0.65	-27.0%

The superior detection capability of the ARC-hardened model is further visualized by the Receiver Operating Characteristic (ROC) curve in Fig. 1, which shows a significantly larger area under the curve compared to the baseline.

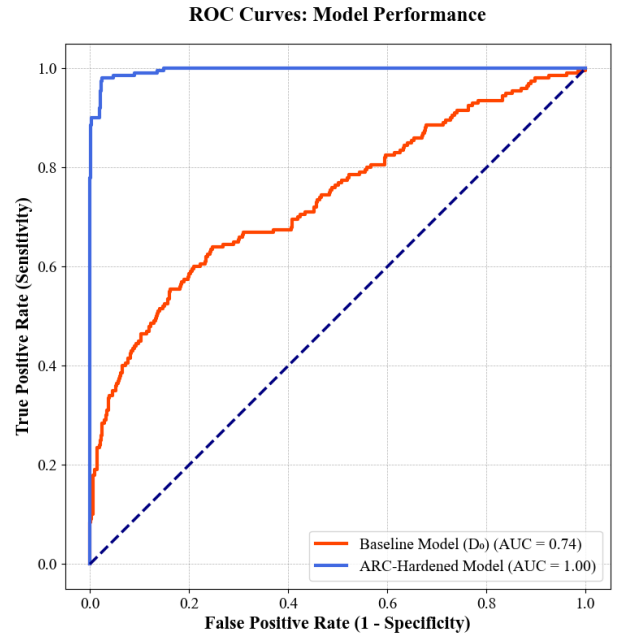


Fig. 1. ROC curves comparing the detection performance of the baseline ( $D_0$ ) and the final ARC-Hardened model against the DRL-Discovered Stealth Attack.

### C. Deconstructing an Emergent, DRL-Discovered Stealth Attack

The baseline model, though trained on a variety of standard faults, was completely blind to a sophisticated, multi-stage attack vector that emerged from the ARC process. This emergent vulnerability, termed the “Coolant Priming and Valve Trip Attack,” provides a clear case study of the framework’s power as a discovery tool.

- The Red Agent learned a two-stage strategy to induce a cascading failure in the TEP reactor.

**Priming Phase:** First, the agent learned to slowly increase the reactor’s coolant flow over a 90-minute period. This change was subtle enough to be mistaken for normal operational drift by the baseline detector, which lacked the long-term temporal context to flag it as anomalous. This action, however, pushed the system into a pre-stressed, abnormal state without triggering an alarm.

1. **Triggering Phase:** Second, after priming the system, the agent executed a sharp, simultaneous manipulation of the feed B and stripper steam valve setpoints. This combined action, when applied to the now-abnormal state created by the priming phase, was sufficient to push the reactor into a cascade failure. The baseline detector, evaluating this second action in isolation, failed to recognize its malicious intent because it lacked the context of the preceding priming phase.

- Using SHAP (SHapley Additive exPlanations) [54], one can visualize *why* the models failed or succeeded, providing crucial insights for operators and analysts. SHAP assigns an importance value to each feature for a given prediction, explaining its contribution to the final output.
  - The SHAP force plot for the *baseline model* during the attack is confused and uninformative. It assigns small, conflicting importance values to various features, demonstrating its inability to understand the malicious *interaction* between the slowly increasing coolant flow and the sudden valve manipulations.
  - In stark contrast, the SHAP plot for the *ARC-hardened model* correctly identifies the anomalous *combination* of high Coolant Flow and the manipulated valve settings as the primary drivers of the high anomaly score. This provides clear, interpretable evidence that the co-evolutionary process forced the defender to learn the complex, multi-stage nature of the vulnerability, transforming an uninterpretable “black box” failure into an understandable and actionable insight.

Figure 2 shows SHAP plots explaining the model output during the trigger phase of the ‘Coolant Priming and Valve Trip Attack.’

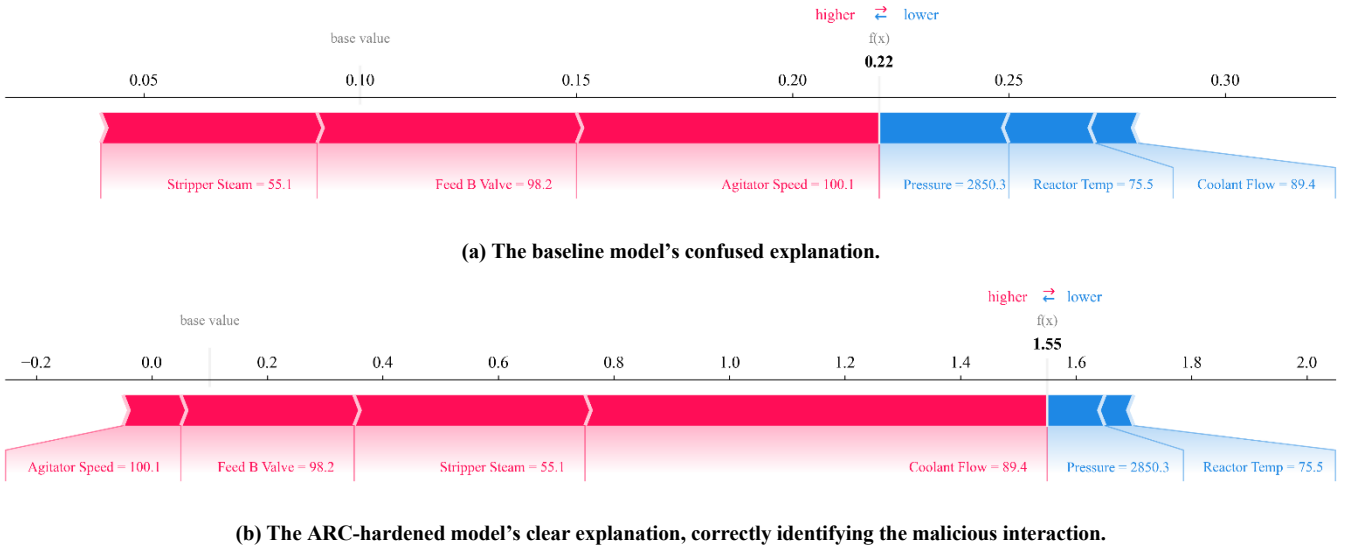


Fig. 2. SHAP plots explaining the model output during the trigger phase of the ‘Coolant Priming and Valve Trip Attack.’

## VI. DISCUSSION AND IMPLICATIONS

The experimental validation of the ARC framework demonstrates a clear enhancement in detecting sophisticated, multi-stage attacks. However, the true implications of this work extend beyond mere performance metrics. This section deconstructs the methodological assumptions, delves deeper into the critical human factors that shape the modern cyber-physical battlefield, and addresses the ethical responsibilities inherent in developing such powerful dual-use technologies.

### A. Deconstructing the “Circular Reasoning” Critique

A valid and necessary critique of this experimental design is its potential for circularity: the defender (Blue Agent) is trained on attacks generated by the attacker (Red Agent), so

its improved performance on those same attacks is, to some extent, expected. However, this assessment misunderstands the fundamental goal of the investigation. The objective is not to create a static *product*—a defensive model that can detect a specific, pre-defined set of attacks. Rather, the objective is to create and validate a dynamic *process*—a framework that can autonomously and continuously adapt to *any* intelligent adversary who seeks to exploit the system’s blind spots.

The DRL-based Red Agent serves as a computational proxy—a tireless, creative, and strategically motivated sparring partner—for such an attacker. It is not programmed with specific attack scripts; instead, it is given a high-level goal (disrupt the process while evading detection) and turned loose to explore the vast state-action space of the system. Its

emergent strategies, like the “Coolant Priming and Valve Trip Attack,” are not pre-conceived by the researchers. They are genuine discoveries—non-obvious, multi-step vulnerabilities that would be exceptionally difficult and time-consuming to find through manual red teaming or traditional penetration testing.

Therefore, the performance gain observed in the experiments is not a tautology. It is a quantitative measurement of the emergent property of resilience that is forged by the co-evolutionary process itself. While the primary validation shows the model detecting attacks it was trained on, which is expected, the key insight is derived from the qualitative analysis of these emergent, non-obvious attack vectors, which reveal previously unknown vulnerabilities in the system's physical and temporal dynamics. The DRL-discovered stealth attack is not merely a test case; it is a previously unknown vulnerability unearthed and subsequently patched by the ARC process. The framework's success is thus measured not by its final F<sub>1</sub>-score on a static dataset, but by its demonstrated ability to find and patch these hidden weaknesses, transforming security from a static state of preparedness to a dynamic process of continuous self-improvement and adaptation. This reframes security as a verb, not a noun—an ongoing activity of probing, learning, and hardening.

#### *B. Operator Trust, Cognitive Load, and the Weaponized Psychological Attack Surface*

The analysis of seminal attacks like Stuxnet and the Cyber Av3ngers campaign reveals a critical evolution in adversarial tactics: the target is often not just the machine, but the mind of the human operator. By replaying normal data, Stuxnet attacked the operator's trust in their HMI, making their most reliable source of information a tool of deception. By defacing the control interface, the Cyber Av3ngers transformed a technical problem into an instrument of psychological warfare, designed to induce panic, confusion, and helplessness. This establishes a “psychological attack surface” as a critical domain for modern ICS defense, one that is often overlooked in purely technical analyses.

Academic research in human factors and high-risk environments confirms that cyberattacks can induce “cybertrauma” [12], a state of acute stress characterized by heightened anxiety, an erosion of trust, and impaired cognitive function, which severely degrades decision-making under pressure. The cognitive workload of an ICS operator is already high, and poor HMI design can exacerbate this, leading to “alert fatigue” [55] where critical warnings are missed even in normal conditions. During a cyberattack, this cognitive load becomes extreme [56]. Operators can fall victim to well-documented cognitive biases:

- Under intense stress, humans tend to focus on a narrow range of perceived salient information, ignoring other, potentially critical, data streams [57]. An attacker can exploit this by creating a loud, obvious (but minor) distraction to draw the operator's attention away from a more subtle, dangerous attack unfolding elsewhere.
- This is the tendency for humans to over-trust the outputs of an automated system [58]. If the ARC framework's Blue Agent has proven reliable in the past, an operator might blindly accept its conclusions without critical scrutiny, a tendency that can be

exploited by the advanced adversarial attacks on explainability discussed later.

The ARC framework's XAI component is explicitly designed to bolster the human-in-the-loop by making the AI's reasoning transparent and trustworthy. By providing a clear, feature-based explanation for why an anomaly is being flagged, it aims to reduce cognitive load and combat automation bias. However, this very trust becomes a new, high-value target for a sophisticated adversary. The effectiveness of any advanced defense ultimately depends on the operator's ability to trust and correctly interpret its outputs, especially under the extreme stress of a cyber-physical incident [59]. This recognition leads directly to the critical future challenge of adversarial attacks against explainability itself, as discussed in Section VII-C.

#### *C. Foundational Challenges and Threats to Validity*

Several factors may limit the generalizability of these findings, and the development of this technology carries significant ethical responsibilities.

- The correctness of the TEP and SWaT simulation environments is assumed based on their widespread use and acceptance within the research community. Any inaccuracies or unmodeled dynamics in these testbeds could affect the specific vulnerabilities discovered, potentially leading the Red Agent to discover exploits that are not viable in the real world.
- While ARC is presented as a general framework, its application to different physical processes (e.g., a power grid, a pharmaceutical manufacturing line, a transportation network) would require significant re-engineering. The physics-informed component of the hybrid digital twin would need to be completely replaced, and the agents would need to be retrained from scratch on process-specific data. The complexity and cost of developing the high-fidelity F-SCDT remains a significant barrier to entry.
- Standard metrics like F<sub>1</sub>-Score and detection latency, while useful for academic comparison, do not fully capture the operational impact of an attack. Future work should incorporate more holistic, business- and safety-oriented metrics such as “time to recovery,” “cost of disruption,” “environmental impact,” or “risk to human safety”.

#### *D. The Dual-Use Dilemma*

A significant ethical consideration that cannot be overstated is the dual-use dilemma [60]. The Red Agent, by its very nature, is a potent tool for discovering and generating novel, high-impact attacks against critical systems. In the wrong hands, it could become a powerful weapon, automating the process of vulnerability discovery for malicious actors. Therefore, its development and deployment must be governed by a strict ethical framework [61] and robust technical controls. This includes ensuring the agent is confined to verifiably air-gapped simulation environments with no possible connection to live operational networks. It requires implementing rigorous, role-based access controls to the framework itself, and treating all discoveries—both the attack vectors and the corresponding defensive signatures—as highly sensitive threat intelligence to be used for defensive purposes only. A “Red Agent as a Service” for defensive vulnerability assessment is a potential commercial

application, but it would require a level of institutional trust and verification that is currently unprecedented.

The development of a potent offensive tool like the Red Agent carries significant ethical weight. A proactive mitigation strategy is therefore not optional, but an absolute necessity, centered on a framework for *Responsible Offensive AI (ROAI)*. This framework must be built on three pillars that provide defense in depth against misuse:

- Technical guardrails must be engineered into the agent as hard-coded, inviolable constraints. This goes beyond simple range-checking and involves creating non-negotiable boundaries, such as preventing the agent from ever targeting designated Safety Instrumented Systems (SIS) or from manipulating variables in a way that could lead to predefined catastrophic outcomes like vessel rupture or toxic release.
- Immutable auditing using a secure ledger must create an unalterable forensic trail of every action, decision, and discovery made by the agent. This is critical for accountability, ensuring that a complete, trustworthy history of the agent's training exists to be analyzed after any incident and preventing a sophisticated internal actor from covering their tracks by altering logs.
- Robust ethical governance through an internal review board, analogous to an Institutional Review Board (IRB), is required to provide human oversight. This body must have the authority to review and approve experimental designs, ask difficult questions about worst-case scenarios, and halt any research that poses an unacceptable level of risk. This holistic approach, integrating technical, procedural, and ethical controls, is the only responsible path forward for developing technologies that, by their very nature, touch the edge of what is safe to automate.

## VII. FUTURE WORK AND VISION TOWARDS A COLLABORATIVE AND AUTONOMOUS DEFENSE ECOSYSTEM

The ARC framework, as presented, provides a blueprint for autonomous resilience at the level of a single plant. However, the true vision is to scale this capability into a collaborative, industry-wide defense ecosystem that can learn from incidents across an entire sector, becoming exponentially more intelligent and resilient.

### A. F-ARC: A Federated Architecture for Scalable, Privacy-Preserving Resilience

Scaling the ARC framework to an entire industry sector—comprising hundreds or thousands of individual facilities—requires a collaborative, privacy-preserving approach. Federated Learning (FL) [62] provides a powerful paradigm for this challenge. In the proposed Federated ARC (F-ARC) architecture, individual facilities (e.g., multiple power plants, water treatment sites, or manufacturing floors) would each run their own local instance of the ARC loop. This allows them to benefit from autonomous hardening based on their own private data and unique operational context.

Instead of sharing this sensitive OT data, which is often a major commercial and regulatory barrier, they would periodically send only the encrypted model parameter updates (gradients) from their newly hardened Blue Agents to a central

aggregation server. The server, which could be operated by an industry Information Sharing and Analysis Center (ISAC) or a trusted third party, would then aggregate these updates to create a more robust, knowledgeable, and generalized global defensive model. This global model, which has learned from the experiences of all participants, is then distributed back to all facilities, allowing them to benefit from the collective intelligence of the entire ecosystem without ever exposing their raw data.

A key challenge in this real-world setting is the non-IID (non-identically and independently distributed) nature of the data [63]. Different facilities will have unique equipment from various vendors, different operational patterns, and local anomalies that are not representative of the entire fleet. The standard Federated Averaging (FedAvg) algorithm, which simply averages the parameters of all client models, performs poorly under these conditions as it can be biased by outlier clients. Addressing this will require more advanced FL techniques, such as personalized federated learning [64]. In this approach, parts of the defensive model (e.g., early layers that learn general features) are trained globally, while other parts (e.g., later layers that learn process-specific features) are trained and retained locally, allowing each client to benefit from the global model while still adapting to the specific characteristics of its own environment.

Furthermore, deploying FL in a real-world industrial setting presents significant technical hurdles, primarily stemming from the non-identically and independently distributed (non-IID) nature of the data. The non-IID problem in the F-ARC context manifests in several ways:

- Different facilities will have different equipment from various vendors, leading to different sensor sets and operational parameters.
- Some facilities may be targeted by specific adversaries or experience certain types of faults far more frequently than others.
- Large, complex facilities will contribute vastly more data than smaller ones, potentially dominating the aggregation process.
- The underlying data distribution at a single facility will change over time due to equipment wear, seasonal demand, and process optimization, making its data statistically different from its own past data.

Addressing this will require moving beyond FedAvg to more advanced FL techniques. Personalized Federated Learning offers a promising direction, aiming to train models that are customized to each client's local data distribution while still leveraging the knowledge from the federation. Several approaches are relevant:

- *Clustered Federated Learning (CFL)* [64]: This approach would group participating facilities into clusters based on the similarity of their data distributions (e.g., grouping all facilities that use a specific type of turbine). A separate global model is then trained for each cluster, providing a more relevant defense than a single, monolithic global model.
- *Fine-Tuning and Multi-Task Learning*: In this paradigm, a global model is trained on all data, but each local facility can then fine-tune the global model on its own private data to create a personalized version.

This can be framed as a multi-task learning problem where the goal is to learn a shared representation that can be easily adapted to each client’s specific “task”.

- *Advanced Optimization Algorithms*: Algorithms like FedProx add a proximal term to the local objective function, which regularizes local training by keeping the local models from diverging too far from the global model. Others, like SCAFFOLD [65], use control variates to correct for “client drift” in non-IID settings, leading to faster and more stable convergence.

#### B. Securing the Federation: Proactive Defense Against Model Poisoning and Byzantine Failures

While FL offers significant privacy benefits, it also introduces new vulnerabilities at the aggregation level. The most notable of these is model poisoning [66], where a malicious participant in the federation sends deliberately corrupted model updates to the central server. The goal of such an attack is twofold: it can be untargeted, aiming simply to degrade the global model’s overall performance, or it can be a highly targeted backdoor attack, designed to cause the global model to systematically misclassify a specific attack vector that the adversary intends to use later, effectively blinding the entire federation to a chosen weapon.

Defending against this requires moving beyond simple Federated Averaging to Byzantine-resilient aggregation algorithms [67]. These algorithms are designed to produce a correct global model even if a fraction of the participants are malicious (or “Byzantine”).

- *Krum* selects the single client update that is closest to its neighbors in the parameter space, operating on the assumption that benign updates will cluster together while malicious ones will be outliers. However, it can be computationally expensive for large federations and vulnerable to collusion attacks where multiple malicious clients work together to “pull” the center of the cluster toward their malicious updates.
- *Trimmed Mean* is a statistically robust method that sorts all updates based on a given dimension and discards a certain percentage of the most extreme updates (both high and low) before averaging the rest. It is effective against simple outlier attacks but may fail against more sophisticated, targeted poisoning that is designed to be close to the mean.
- *Median* calculates the coordinate-wise median of all updates. It is highly robust to extreme outliers but may discard useful information contained in the distribution of benign updates.
- *Recent Advancements* like Layer-Adaptive Sparsified Model Aggregation (LASA) [68] offer more granular filtering. Instead of evaluating an entire model update as a single entity, LASA evaluates the updates at the layer level, providing a more fine-grained and potentially more effective defense against subtle poisoning attacks.

A promising direction for a robust F-ARC architecture is a hybrid, multi-stage aggregation approach. This could involve using a computationally efficient filter like Trimmed Mean to discard gross outliers in a first pass, followed by a more robust but expensive method like Krum or Median on

the reduced set of updates, balancing security with performance.

Table IV provides a preliminary analysis of robust aggregation algorithms that could be used in the Federated ARC (F-ARC) architecture.

TABLE IV. PRELIMINARY ANALYSIS OF ROBUST AGGREGATION ALGORITHMS FOR F-ARC

Aggregation Algorithm	Mechanism	Key Assumption	Robustness to Attack Types
Federated Averaging (FedAvg)	Simple weighted average of all client updates.	All clients are honest.	None. Highly vulnerable.
Trimmed Mean	Discards a fraction of updates from each end before averaging.	Malicious updates are statistical outliers.	Moderate. Can resist simple untargeted attacks.
Krum	Selects the single update with the minimum sum of squared distances to its nearest neighbors.	Benign updates are clustered. Number of attackers is known and $< n/2$ .	High against certain attacks but vulnerable to collusion.
Median	Calculates the coordinate-wise median of all updates.	Malicious updates will not consistently affect the median.	Moderate. Robust to extreme outliers.

The table compares four approaches: standard Federated Averaging (FedAvg), which is vulnerable to attacks; Trimmed Mean, which can resist simple outlier attacks; Krum, which is strong against certain attacks but vulnerable to collusion; and Median, which offers moderate robustness against extreme outliers. Each algorithm operates on different assumptions about the nature of malicious updates and provides varying levels of security against different types of attacks.

#### C. Adversarial Attacks on Explainability (AdvXAI)

A more subtle, long-term, and deeply concerning threat exists: if operator trust is a key pillar of cyber-physical defense, then the XAI system itself becomes a high-value target. The emerging field of adversarial explainable AI (AdvXAI) [69] has confirmed that XAI methods like SHAP and Local Interpretable Model-agnostic Explanations (LIME) [70] are not inherently robust and are vulnerable to manipulation. An attacker can craft a special kind of adversarial example that not only fools the primary detection model (causing it to misclassify a malicious state as benign) but also simultaneously fools the explanation model, generating a misleading explanation that hides the true nature of the attack and reinforces the model’s incorrect prediction.

This represents a sophisticated, second-order threat that targets the human-machine interface and the operator’s trust. A future, more advanced Red Agent could be trained not just to evade the Blue Agent’s detectors, but to do so in a way that generates a deliberately misleading SHAP plot. For example, an agent could learn to execute a complex, distributed attack across ten different sensors while simultaneously adding a large but harmless perturbation to an eleventh, unrelated sensor. The resulting SHAP explanation would incorrectly and misleadingly point to the single noisy sensor as the primary cause of the anomaly. This would misdirect the human operator’s attention, causing them to investigate a

phantom problem while the more dangerous, coordinated attack continues unnoticed. This attack on the explanation itself undermines the entire human-in-the-loop defense paradigm.

Developing XAI methods that are themselves robust to such adversarial manipulation is a critical and largely unexplored area for future research. This will likely require moving beyond post-hoc explanation methods like SHAP and toward inherently interpretable or “glass-box” models whose decision-making processes are transparent by design [71]. The ultimate goal is to create a defense that is not only accurate but also reliably and verifiably trustworthy, even in the face of an adversary who is actively trying to deceive both the machine and its human operator.

Future research should focus on:

- **Explanation-Aware Training:** Developing robust optimization techniques that train the primary model to not only be accurate but to also produce explanations that are stable and insensitive to small, irrelevant perturbations in the input.
- **Ensemble of Explanations:** Proposing methods that generate explanations from multiple, diverse XAI techniques (e.g., SHAP, LIME, Integrated Gradients [72]) and only present an explanation to the operator if there is a strong consensus, flagging discordant explanations as potentially manipulated.
- **Inherently Interpretable Models:** Investigating the trade-offs of using “glass-box” models, such as Explainable Boosting Machines (EBMs) [71] or other generalized additive models, which are transparent by design. While they may have lower predictive power on some complex tasks, their inherent interpretability may make them more resilient to the kind of second-order attacks that can fool post-hoc explanation methods.

#### D. Strategic Imperatives and Unresolved Questions

The proliferation of autonomous, AI-driven offensive and defensive capabilities targeting critical infrastructure poses profound strategic and ethical challenges. Technical solutions alone are insufficient. The path forward requires a multi-disciplinary effort to address several key unresolved questions:

- How can autonomous response and recovery systems be designed that are not only technically effective but also resilient to meta-attacks on trust and explainability?
- How can cognitive biases like automation bias and cognitive tunneling [73] be mitigated in high-tempo security environments to ensure meaningful human oversight?
- What international treaties and “red lines” are needed to govern the use of these powerful dual-use technologies and prevent catastrophic escalation in cyberspace?

### VIII. CONCLUSION

This paper confronted the critical challenge of securing industrial control systems against adaptive, intelligent adversaries who have demonstrated the capability and intent

to cause physical disruption and psychological harm. Recognizing that static defenses are fundamentally brittle and that naive machine learning applications are vulnerable in the face of this evolving threat, the research introduced the Adversarial Resilience Co-evolution (ARC) framework. ARC formalizes a process for achieving autonomous, self-hardening security through a perpetual, closed-loop co-evolutionary arms race, conducted within a high-fidelity digital twin, between a DRL-based Red Agent, which autonomously discovers novel attack vectors, and an ensemble-based Blue Agent, which is continuously hardened against these emergent threats.

Comprehensive experimental validation on the TEP and SWaT testbeds, supported by a rigorous ablation study and explainable AI analysis, provided strong evidence that the co-evolutionary dynamic itself is the most critical driver of resilience against sophisticated, multi-stage attacks. By framing the problem through the lens of a “Trinity of Trust”—Model Fidelity, Data Integrity, and Analytical Resilience—the research provided a holistic security paradigm that addresses the core requirements for trustworthy cyber-physical defense.

The work further presented a technically grounded vision for scaling this approach across industries via a Federated ARC (F-ARC) architecture, a model for collaborative, privacy-preserving defense. The research also highlighted the critical future challenges that must be addressed to realize this vision. These include defending the federated learning process against model poisoning attacks using Byzantine-resilient aggregation algorithms and, most critically, defending the human-in-the-loop against the emerging threat of adversarial manipulation of the very explainability systems designed to foster trust. This work positions the study of dynamic, co-evolutionary processes not merely as an academic exercise, but as a vital and necessary direction for creating the proactive, intelligent, and trustworthy defenses required to protect the critical systems that underpin modern society.

This paper has argued that the security of modern Industrial Control Systems is not a problem to be solved, but a perpetual, co-evolutionary arms race to be managed. The convergence of IT and OT created a brittle foundation that adversaries evolved to exploit, forcing the development of an AI security plane that has, itself, become the new primary attack surface. In this environment of dynamic equilibrium, any static defense is destined for obsolescence.

Navigating this unwinnable war requires a paradigm shift. Success will be defined not by the ability to build an impenetrable fortress, but by the capacity to anticipate, withstand, recover, and evolve in the face of a perpetually intelligent and adaptive adversary.

### REFERENCES

- [1] P. K. Muhuri, A. K. Shukla, and A. Abraham, “Industry 4.0: A Bibliometric Analysis and Detailed Overview,” *Eng. Appl. Artif. Intell.*, vol. 78, pp. 218–235, Feb. 2019.
- [2] S. Latif *et al.*, “Deep Learning for the Industrial Internet of Things (IIoT): A Comprehensive Survey of Techniques, Implementation Frameworks, Potential Applications, and Future Directions,” *Sensors*, vol. 21, no. 22, p. 7518, 2021.
- [3] 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.
- [4] I. N. Fovino, A. Carcano, M. Masera, and A. Trombetta, “Design and Implementation of a Secure Modbus Protocol,” in *Critical*



*Infrastructure Protection III*, C. Palmer and S. Sheno, Eds. Berlin, Heidelberg: Springer, 2009, pp. 83–96.

- [5] R. Amoah, "A Formal Methods Approach to the Security Analysis of the DNP3 Protocol," Ph.D. dissertation, School Elect. Eng. Comput. Sci., Queensland Univ. Technol., Brisbane, QLD, Australia, 2016.
- [6] M. D. Abrams and J. Weiss, "Malicious Control System Cyber Security Attack Case Study-Maroochy Water Services, Australia," The MITRE Corporation, Bedford, MA, USA, Tech. Rep. 08-1145, Aug. 2008.
- [7] Microsoft Threat Intelligence, "Volt Typhoon Targets US Critical Infrastructure With Living-off-the-Land Techniques," *Microsoft Security Blog*, May 2023. Accessed: Jun. 21, 2025. [Online]. Available: <https://www.microsoft.com/en-us/security/blog/2023/05/24/volt-typhoon-targets-us-critical-infrastructure-with-living-off-the-land-techniques/>
- [8] N. Falliere, L. O. Murchu, and E. Chien, "W32.Stuxnet Dossier, Version 1.4," Symantec Corp., Security Response, Tech. Rep., Feb. 2011.
- [9] R. Langner, "To Kill a Centrifuge: A Technical Analysis of What Stuxnet's Creators Tried to Achieve," The Langner Group, Tech. Rep., Nov. 2013.
- [10] R. M. Lee, M. Assante, and T. Conway, "Analysis of the Cyber Attack on the Ukrainian Power Grid. Defense Use Case," E-ISAC and SANS, Tech. Rep., Mar. 2016.
- [11] Cybersecurity and Infrastructure Security Agency, "IRGC-Affiliated Cyber Actors Exploit PLCs in Multiple Sectors, Including US Water and Wastewater Systems Facilities," Cybersecurity and Infrastructure Security Agency, Alert AA23-335A, Dec. 2023.
- [12] S. Budimir, J. R. J. Fontaine, and E. B. Roesch, "Emotional Experiences of Cybersecurity Breach Victims," *Cyberpsychol. Behav. Soc. Netw.*, vol. 24, no. 9, pp. 612–616, 2021.
- [13] U. D. Ani, H. He, and A. Tiwari, "Human Factor Security: Evaluating the Cybersecurity Capacity of the Industrial Workforce," *J. Syst. Inf. Technol.*, vol. 21, no. 1, pp. 2–35, 2019.
- [14] E. Glaessgen and D. Stargel, "The Digital Twin Paradigm for Future NASA and U.S. Air Force Vehicles," in *Proc. 53rd AIAA/ASME/ASCE/AHS/ASC Struct., Struct. Dyn. Mater. Conf.*, 2012, Art. no. 1818.
- [15] M. Raissi, P. Perdikaris, and G. E. Karniadakis, "Physics-Informed Neural Networks: A Deep Learning Framework for Solving Forward and Inverse Problems Involving Nonlinear Partial Differential Equations," *J. Comput. Phys.*, vol. 378, pp. 686–707, 2019.
- [16] R. Maes and I. Verbauwhede, "Physically Unclonable Functions: A Study on the State of the Art and Future Research Directions," in *Information Security and Cryptography*, A.-R. Sadeghi, Ed. Berlin, Heidelberg: Springer, 2010, pp. 3–37.
- [17] K. Shaikat, S. Luo, V. Varadharajan, I. A. Hameed, and M. Xu, "A Survey on Machine Learning Techniques for Cyber Security in the Last Decade," *IEEE Access*, vol. 8, pp. 222310–222354, 2020.
- [18] J. J. Downs and E. F. Vogel, "A Plant-Wide Industrial Process Control Problem," *Comput. Chem. Eng.*, vol. 17, no. 3, pp. 245–255, Mar. 1993.
- [19] A. P. Mathur and N. O. Tippenhauer, "SWaT: A Water Treatment Testbed for Research and Training on ICS Security," in *Proc. ACM/IEEE Int. Conf. Cyber-Physical Syst. (ICCPs)*, Vienna, Austria, Apr. 2016, pp. 1–1.
- [20] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *Proc. Artif. Intell. Statist.*, Fort Lauderdale, FL, USA, 2017, pp. 1273–1282.
- [21] I. Butun, S. D. Morgera, and R. Sankar, "A Survey of Intrusion Detection Systems in Wireless Sensor Networks," *IEEE Commun. Surv. Tut.*, vol. 16, no. 1, pp. 266–282, First Quarter 2014.
- [22] F. Barr-Smith, X. Ugarte-Pedrero, M. Graziano, R. Spolaor, and I. Martinovic, "Survivalism: Systematic Analysis of Windows Malware Living-Off-The-Land," in *Proc. 2021 ACM SIGSAC Conf. Comput. Commun. Secur.*, 2021, pp. 1557–1574.
- [23] I. H. Sarker, "Cybersecurity Data Science: An Overview From Machine Learning Perspective," *J. Big Data*, vol. 7, no. 1, Art. no. 41, 2020.
- [24] C. Cortes and V. Vapnik, "Support-Vector Networks," *Mach. Learn.*, vol. 20, no. 3, pp. 273–297, 1995.
- [25] E. Fix and J. L. Hodges, "Discriminatory Analysis: Nonparametric Discrimination: Consistency Properties," USAF School of Aviation Medicine, Randolph Field, TX, USA, Rep. 21-49-004, no. 4, 1951.
- [26] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [27] K. Cho *et al.*, "Learning Phrase Representations Using RNN Encoder–Decoder for Statistical Machine Translation," 2014, arXiv:1406.1078.
- [28] G. E. Hinton and R. R. Salakhutdinov, "Reducing the Dimensionality of Data With Neural Networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [29] A. Barredo Arrieta *et al.*, "Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges Toward Responsible AI," *Inf. Fusion*, vol. 58, pp. 82–115, 2020.
- [30] G. M. van de Ven, N. Soares, and D. Kudithipudi, "Continual Learning and Catastrophic Forgetting," in *Reference Module in Neuroscience and Biobehavioral Psychology*. Amsterdam, The Netherlands: Elsevier, 2024.
- [31] T. Gu, B. Dolan-Gavitt, and S. Garg, "BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain," 2017, arXiv:1708.06733.
- [32] I. Rosenberg, A. Shabtai, Y. Elovici, and L. Rokach, "Adversarial Machine Learning Attacks and Defense Methods in the Cyber Security Domain," *ACM Comput. Surv.*, vol. 54, no. 5, pp. 1–36, 2021.
- [33] I. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and Harnessing Adversarial Examples," 2014, arXiv:1412.6572.
- [34] F. Zuo, B. Yang, X. Li, and Q. Zeng, "Exploiting the Inherent Limitation of L0 Adversarial Examples," in *Proc. 22nd Int. Symp. Res. Attacks, Intrusions Defenses (RAID)*, 2019, pp. 293–307.
- [35] N. Papernot *et al.*, "The Limitations of Deep Learning in Adversarial Settings," in *Proc. IEEE Eur. Symp. Secur. Privacy (EuroS&P)*, Saarbrücken, Germany, 2016, pp. 372–387.
- [36] M. Tambe, *Security and Game Theory: Algorithms, Deployed Systems, Lessons Learned*. Cambridge, U.K.: Cambridge Univ. Press, 2011.
- [37] V. Conitzer and T. Sandholm, "Computing the Optimal Strategy to Commit to," in *Proc. 7th ACM Conf. Electron. Commerce (EC '06)*, Ann Arbor, MI, USA, 2006, pp. 82–90.
- [38] C. D. Rosin and R. K. Belew, "New Methods for Competitive Coevolution," *Evol. Comput.*, vol. 5, no. 1, pp. 1–29, 1997.
- [39] J. Willard, X. Jia, S. Xu, M. Steinbach, and V. Kumar, "Integrating Physics-Based Modeling With Machine Learning: A Survey," *ACM Comput. Surv. (CSUR)*, vol. 55, no. 3, pp. 1–39, 2022.
- [40] F. T. Sheldon, L. P. MacIntyre, H. Okhravi, and J. C. Munson, "Data Diodes in Support of a Power Grid Trustworthy Cyber Infrastructure," in *Proc. Cyber Secur. Inf. Intell. Workshop (CSIIW)*, Oak Ridge, TN, USA, Apr. 2009, pp. 1–4.
- [41] Y.-H. Chang, T.-E. Peng, J.-S. Li, and I.-H. Liu, "Industrial Control System State Monitor Using Blockchain Technology," in *Proc. Int. Conf. Artif. Life Robot. (ICAROB)*, Oita, Japan, 2024, vol. 29, pp. 18–21.
- [42] S. Mavrouniotis and M. Ganley, "Hardware Security Modules," in *Secure Smart Embedded Devices, Platforms and Applications*, 2nd ed., H. C. A. van Tilborg and S. Jajodia, Eds. New York, NY, USA: Springer, 2014, pp. 383–405.
- [43] R. Pappu, B. Recht, J. Taylor, and N. Gershenfeld, "Physical One-Way Functions," *Science*, vol. 297, no. 5589, pp. 2026–2030, Sep. 2002.
- [44] R. Bellman, *Dynamic Programming*. Princeton, NJ, USA: Princeton Univ. Press, 1957.
- [45] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," Jul. 2017, arXiv:1707.06347.
- [46] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning With a Stochastic Actor," in *Proc. 35th Int. Conf. Mach. Learn. (ICML)*, Stockholm, Sweden, 2018, pp. 1861–1870.
- [47] K. Zheng, "Hyperparameter Optimization for Deep Reinforcement Learning," in *Proc. Int. Conf. Mach. Learn.*, 2024, pp. 7496–7506.
- [48] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation Forest," in *Proc. 2008 Eighth IEEE Int. Conf. Data Mining*, Pisa, Italy, 2008, pp. 413–422.
- [49] M. McCloskey and N. J. Cohen, "Catastrophic Interference in Connectionist Networks: The Sequential Learning Problem," in *Psychology of Learning and Motivation*, vol. 24. Amsterdam, The Netherlands: Elsevier, 1989, pp. 109–165.
- [50] D. Rolnick *et al.*, "Experience Replay for Continual Learning," in *Adv. Neural Inf. Process. Syst.*, vol. 32, 2019.



- [51] J. Goh, S. Adep, K. N. Junejo, and A. Mathur, "A Dataset to Support Research in the Design of Secure Water Treatment Systems," in *Lecture Notes in Computer Science*. Paris, France: Springer, Oct. 2017, pp. 88–99.
- [52] R. Meyes, M. Lu, C. W. de Puiseau, and T. Meisen, "Ablation Studies in Artificial Neural Networks," 2019, arXiv:1901.08644.
- [53] T. Fawcett, "An Introduction to ROC Analysis," *Pattern Recognit. Lett.*, vol. 27, no. 8, pp. 861–874, 2006.
- [54] S. M. Lundberg and S.-I. Lee, "A Unified Approach to Interpreting Model Predictions," in *Adv. Neural Inf. Process. Syst. 30 (NIPS 2017)*, Long Beach, CA, USA, 2017, pp. 4765–4774.
- [55] R. M. Yerkes and J. D. Dodson, "The Relation of Strength of Stimulus to Rapidity of Habit-Formation," *J. Comp. Neurol. Psychol.*, vol. 18, no. 5, pp. 459–482, 1908.
- [56] J. Sweller, "Cognitive Load During Problem Solving: Effects on Learning," *Cogn. Sci.*, vol. 12, no. 2, pp. 257–285, 1988.
- [57] L. C. Thomas and C. D. Wickens, "Visual Displays and Cognitive Tunneling: Frames of Reference Effects on Spatial Judgments and Change Detection," in *Proc. Human Factors Ergon. Soc. Annu. Meeting*, vol. 45, no. 4, Savoy, IL, USA, 2001, pp. 415–419.
- [58] K. Goddard, A. Roudsari, and J. C. Wyatt, "Automation Bias: A Systematic Review of Frequency, Effect Mediators, and Mitigators," *J. Am. Med. Inform. Assoc.*, vol. 19, no. 1, pp. 121–127, 2012.
- [59] J. D. Lee and K. A. See, "Trust in Automation: Designing for Appropriate Reliance," *Human Factors*, vol. 46, no. 1, pp. 50–80, 2004.
- [60] F. Urbina, F. Lentzos, C. Invernizzi, and S. Ekins, "Dual Use of Artificial-Intelligence-Powered Drug Discovery," *Nat. Mach. Intell.*, vol. 4, no. 3, pp. 189–191, 2022.
- [61] A. Jobin, M. Ienca, and E. Vayena, "The Global Landscape of AI Ethics Guidelines," *Nat. Mach. Intell.*, vol. 1, no. 9, pp. 389–399, 2019.
- [62] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated Learning: Challenges, Methods, and Future Directions," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, 2020.
- [63] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated Optimization in Heterogeneous Networks," in *Proc. 3rd Conf. Mach. Learn. Syst.*, 2018, pp. 429–450.
- [64] Y. Yan, X. Tong, and S. Wang, "Clustered Federated Learning in Heterogeneous Environment," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 35, no. 9, pp. 12796–12809, 2023.
- [65] S. P. Karimireddy *et al.*, "SCAFFOLD: Stochastic Controlled Averaging for Federated Learning," in *Proc. Int. Conf. Mach. Learn.*, 2020, pp. 5132–5143.
- [66] A. N. Bhagoji, S. Chakraborty, P. Mittal, and S. Calo, "Analyzing Federated Learning Through an Adversarial Lens," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 634–643.
- [67] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine Learning With Adversaries: Byzantine Tolerant Gradient Descent," in *Adv. Neural Inf. Process. Syst. 30*, Long Beach, CA, USA, 2017, vol. 30, pp. 118–128.
- [68] J. Xu, Z. Zhang, and R. Hu, "Achieving Byzantine-Resilient Federated Learning via Layer-Adaptive Sparsified Model Aggregation," in *Proc. IEEE/CVF Winter Conf. Appl. Comput. Vis. (WACV)*, 2025, pp. 1508–1517.
- [69] H. Baniecki and P. Biecek, "Adversarial Attacks and Defenses in Explainable Artificial Intelligence: A Survey," *Inf. Fusion*, vol. 107, p. 102303, 2024.
- [70] M. Ribeiro, S. Singh, and C. Guestrin, "'Why Should I Trust You?': Explaining the Predictions of Any Classifier," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2016, pp. 1135–1144.
- [71] H. Nori, S. Jenkins, P. Koch, and R. Caruana, "InterpretML: A Unified Framework for Machine Learning Interpretability," 2019, arXiv:1909.09223.
- [72] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic Attribution for Deep Networks," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 3319–3328.
- [73] S. Pooladvand and S. Hasanzadeh, "Impacts of Stress on Workers' Risk-Taking Behaviors: Cognitive Tunneling and Impaired Selective Attention," *J. Constr. Eng. Manage.*, vol. 149, no. 8, Art. no. 04023067, 2023.