# Scenario-Agnostic Zero-Trust Defense with Explainable Threshold Policy: A Meta-Learning Approach

Yunfei Ge*, Tao Li*, and Quanyan Zhu

*Abstract*—The increasing connectivity and intricate remote access environment have made traditional perimeter-based network defense vulnerable. Zero trust becomes a promising approach to provide defense policies based on agent-centric trust evaluation. However, the limited observations of the agent's trace bring information asymmetry in the decision-making. To facilitate the human understanding of the policy and the technology adoption, one needs to create a zero-trust defense that is explainable to humans and adaptable to different attack scenarios. To this end, we propose a scenario-agnostic zero-trust defense based on Partially Observable Markov Decision Processes (POMDP) and first-order Meta-Learning using only a handful of sample scenarios. The framework leads to an explainable and generalizable trust-threshold defense policy. To address the distribution shift between empirical security datasets and reality, we extend the model to a robust zero-trust defense minimizing the worst-case loss. We use case studies and real-world attacks to corroborate the results.

*Index Terms*—Zero-trust security, meta learning, scenario-agnostic, threshold policy

Fig. 1: *Illustration of the scenario-agnostic zero-trust defense. The meta policy acquired in the meta-learning phase can respond to unknown attack scenarios with fine-tuned adaptation.*

## I. INTRODUCTION

The recent advances in cloud services, data communication, and automation technologies have increased flexibility and efficiency in modern network systems [1]. However, the adoption of smart devices and the Internet of Things (IoT) has brought up new and expanding cyber risks, not just capable of impacting a particular device but creating severe concerns in the whole system [2]. The increasing connectivity, heterogeneity, and dynamic accessing environments inevitably enlarge the attack surface and lead to multiple vulnerabilities that attackers can exploit. In response to the vulnerabilities in traditional perimeter-based network security, modern networks must transform from static and perimeter-based defenses to a zero-trust security framework that forfeits the assumption that everything behind the security perimeter is safe. It eliminates implicit trust in each agent and makes defense decisions based on continuous trust evaluation [3].

However, designing the zero-trust defense is not trivial, and several challenges arise. First, the limited observations of the agent's trace bring information asymmetry in the decision-making. Hence, a quantitative metric measuring the agents' trustworthiness using partial observations is indispensable. Moreover, the zero-trust policy must be generalizable to a family of scenarios. Besides, it is necessary to create a zero-trust defense that is explainable to human operators who develop security solutions based on the reasoning of the defense policy. In this way, explainable and adaptable zero-trust defenses make the configuration of the large-scale network easier and facilitate the broad adoption of the technology.

To equip the zero-trust defense with adaptability under information asymmetry, we propose a new zero-trust defense
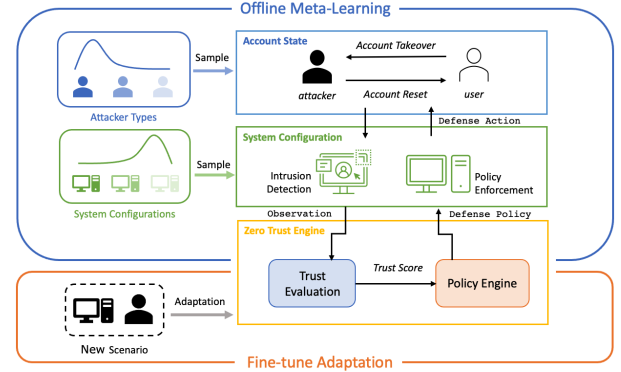
The authors are with Department of Electrical and Computer Engineering, New York University, NY; E-mail: {yg2047,taoli,qz494}@nyu.edu. *Yunfei Ge and Tao Li contributed equally to this work.

taking a threshold form based on the concept of meta-learning [4]. Unlike classical learning paradigms (i.e., reinforcement learning), meta-learning aims to learn a learning strategy using previous training data, rather than a simple decision-making model. In the face of a new task unseen in the training phase, the obtained learning strategy enables the defender to learn a satisfying policy using far less data than from scratch.

Taking inspiration from meta-learning, the proposed ZTD, referred to as scenario-agnostic zero-trust defense (SA-ZTD), enables the defender to adapt to new scenarios (system configurations/attack capabilities) with a modest amount of partial observations using a learning strategy learned from the training experience. To be more specific, we first formulate a zero-trust network security model using parameterized Partially Observed Markov Decision Processes (POMDP) [5], where the parameters represent attack scenarios with distinct system vulnerabilities and the attacker's capabilities.

Since real-world applications involve a large (possibly infinite) number of attack scenarios, it is intractable to compute the optimal policy for each scenario. The proposed SA-ZTD resolves this issue by learning a meta policy and an adaptation mapping (i.e., the learning strategy) using a handful of known scenarios. When deployed in a new attack scenario, the learned mapping can quickly adapt the meta policy to the current environment using a few partial observations. We use the word "agnostic" (whose root means "not known") to emphasize that the adaptation ability is acquired without knowledge of every scenario. An illustration of SA-ZTD is presented in Figure 1. **Our contributions** are as follows. 1) We propose scenario-agnostic zero-trust defense, a generalizable defense strategy that does require access to every scenario. In addition, SA-ZTD leads to explainable trust-threshold policies. 2) To address the distribution shift between empirical security datasets and reality, we propose a robust zero-trust defense based on SA-ZTD that minimizes worst-case loss. 3) A first-order meta-learning algorithm is developed to learn SA-ZTD

efficiently using only a handful of sample scenarios.

**Related Works**: Zero-trust defense has become an emerging trend for addressing the challenges in modern network security. A conceptual zero trust strategy is proposed in [6] to establish the trust notion in cloud computing. [7] investigates a zero-trust architecture for 5G networks utilizing artificial intelligence to provide information security. Authors in [8] have combined zero trust and block-chain to manage IoT device security. Despite the extensive applications, the aforementioned frameworks are conceptual, and the proposed zero-trust defenses highly depend on the underlying system configuration. There is a lack of adaptation/generalization ability in their defense strategy design, and our work is among the first endeavor to investigate adaptable zero-trust defense leveraging meta-learning.

## II. ZERO-TRUST DEFENSE UNDER ASYMMETRIC INFORMATION

Consider account take-over attacks (ATA), where the attacker attempts to compromise the legitimate account using social engineering [9] and zero-day vulnerabilities [10] at every time step. Meanwhile, the system administrator (the defender) can turn the adversarial account into a legitimate one by removing the stored credentials and resetting the account.

The true nature of the account, while revealed to the attacker, remains hidden from the defender: only the attacker knows whether the attack succeeds or not. In contrast, the defender can only observe the footprints of the agent through system alerts and monitoring information, such as intrusion-detection systems. Note that the system alerts do not equate to the actual state of the account, as the monitoring mechanism may produce false positive/negative results.

The asymmetric information structure in ATA, as in other security problems [11], complicates the defender's decision-making process. The defender needs to dynamically evaluate the trust using partial observations and then decide whether disable the account at the policy enforcement point. The above zero-trust defense problem can be formulated as a POMDP defined by the tuple parameterized by $\theta \in \Theta$: $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, T_\theta, O, C, \rho \rangle$ discussed in the following. $k \in \{0, 1, \dots\}$ denotes the discrete time index.

- $\Theta$ denotes the set of all attack **scenarios**, and each parameter $\theta \in \Theta$ captures the system vulnerabilities and attacker capabilities that affect the system transition to be defined later.
- $\mathcal{S} = \{0, 1\}$ denotes the set of **account states**. $s = 0$ stands for the adversarial, while $s = 1$ for the legitimate;
- $\mathcal{A} = \{0, 1\}$ is the set of **defense actions**. $a = 0$ means removing the account credentials and resetting the account. $a = 1$ indicates that the defender takes no action;
- $\mathcal{O} = \{0, 1\}$ denotes the set of **system observations**. $o = 0$ indicates the defender observes a security alert about the anomaly behaviors of the account, and no alert is observed when $0 = 1$;
- $T^\theta(a) \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, for $a \in \mathcal{A}$ is the **system transition matrix** under action $a$. The $ij$-entry of the ma-

trix given indicates the probability that the account of state $i$ being changed to state $j$, i.e., $T_{ij}^\theta(a) = \mathbb{P}\left(s^{k+1} = j | s^k = i, a^k = a, \theta\right), i, j \in \mathcal{S}$.

Since the defense action is binary, we divide the presentation of the transition matrices into two cases below.

1) **Active Defense Case**: $a = 0$. The defender decides to reset the account, and as a result, the attacker loses its foothold within the system with probability $1 - p_a^d$. $p_a^d \in [0, 1]$ is the probability that the attacker bypasses the defense when controlling the account. On the other hand, if the account is legitimate, the attacker can compromise it with probability $p_u^d \in [0, 1]$. In summary, $p_a^d$ and $p_u^d$ represent the attacker's capabilities in launching and sustaining ATA, and the transition matrix under active defense is given by $T_{00}(0) = p_a^d, T_{01}(0) = 1 - p_a^d, T_{10}(0) = p_u^d, T_{11}(0) = 1 - p_u^d$.

2) **Normal Operation Case**: $a = 1$. In this case, no defense action is taken, and $p_u^n \in [0, 1]$ is the probability that the legitimate account is taken over by the attacker. The magnitude of this value reflects the system vulnerability: the larger $p_u^n$, the more vulnerable the system is. Let $p_a^n \in [0, 1]$ be the probability that the attacker in the system without being detected during normal operation (stealthiness), and the transitions are as below: $T_{00}(1) = p_a^n, T_{01}(1) = 1 - p_a^n, T_{10}(1) = p_u^n, T_{11}(1) = 1 - p_u^n$.

The system transition matrices are jointly determined by the system vulnerability ($p_u^n$) and the attacker's capability/stealthiness ($p_a^d, p_u^d, p_a^n$). Hence, each attack scenario is fully captured by the concatenation of the relevant parameters in the transition, i.e., $\theta = (p_a^d, p_u^d, p_a^n, p_u^n) \in [0, 1]^4$.

- $O \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{O}|}$ is the observation matrix whose $ij$-entry is defined as $O_{00} = q_a, O_{01} = 1 - q_a, O_{10} = q_u, O_{11} = 1 - q_u$, where $q_a, q_u \in [0, 1]$ are the detection rate and false alarm rate of the intrusion detection system, respectively.
- $C(s, a) \in \mathbb{R}$ is the defender's cost when implementing $a$ at state $s$. In particular, the cost against the attacker during normal operation $C(1, 0)$ indicates the potential damage of the attack. The defense cost $C(\cdot, 1)$ captures the time delay and performance degradation due to the account reset. The defense performance is evaluated by the cumulative cost discounted by the factor $\rho \in (0, 1)$. The defender's objective is to find an optimal policy that balances security and system performance across different scenarios to be detailed in the next section.

## III. SCENARIO-AGNOSTIC ZERO-TRUST DEFENSE

For a specific attack scenario $\theta$, a POMDP-based zero-trust defense under asymmetric information is developed in [12], where the defender dynamically evaluates the account's trustworthiness using the *trust score (TS)*, defined as the belief that the user is legitimate, i.e., $TS^k := b^k(s = 1)$. The defender's belief is updated in the following Bayesian manner: for $s' \in \mathcal{S}$,

$$b^k(s') = \frac{O_{s', o^{k+1}}(a^k) \sum_{s \in \mathcal{S}} T_{s, s'}(a^k) b^k(s)}{\sum_{s' \in \mathcal{S}} O_{s', o^{k+1}}(a^k) \sum_{s \in \mathcal{S}} T_{s, s'}(a^k) b^k(s)}.$$

Based on this trust evaluation, the defender searches for a defense policy in $\Pi := \{\pi | \pi : \Delta(\mathcal{S}) \to \mathcal{A}\}$ to minimize the expected cumulative cost $U_\theta(\pi) = \mathbb{E}_{a^k \sim \pi(b^k), T^\theta, O}[\sum_{k=0}^\infty \rho^k c(s^k, a^k) \mid b^0]$, where $b^0 \in \Delta(\mathcal{S})$ specifies the initial trust of the account.

A stochastic gradient descent algorithm is proposed in [12] to minimize the cost $U_\theta(\pi)$, leading to a simple approach to zero-trust defense (ZTD). However, the resulting optimal policy is generally scenario-dependent: the defense policy is designed for some specific system configuration and attacker capabilities. The insufficiency of this approach is the lack of **adaptation/generalization** to different scenarios. As observed in our experiments, a slight change in $\theta$ leads to differences in defense policies, meaning that the defense policy for one scenario does not generalize to another.

To equip the ZTD with adaptation/generalization ability, we propose a new zero-trust approach based on the meta-learning idea. As a learning-to-learn approach, meta learning [4] intends to build an internal representation of the policy (meta policy) that is broadly suitable for a collection of different but related scenarios. When deployed in a specific scenario possibly unseen in the training phase, the meta policy can quickly adapt to the current environment using a few data, where the adaptation strategy is learned from past training data. In short, two pillars of meta learning are the meta policy $\pi_{meta}$ and the adaptation mapping $\phi : \Pi \times \Theta \to \Pi$, respectively. Leveraging the above notions, we present a meta-learning-based ZTD in the following.

**Definition 1** (Scenario-Agnostic Zero-Trust Defense). *A pair $\langle \pi_{meta}, \phi \rangle$ is said to be a scenario-agnostic zero-trust defense (SA-ZTD) with respect to a scenario distribution $p \in \Delta(\Theta)$ if the pair solves for the minimization problem below*

$$\min_{\pi, \phi} \mathbb{E}_{\theta \sim p}[U_\theta(\phi(\pi, \theta))] \tag{1}$$

Such a defense is scenario-agnostic in that solving for (1) (approximately) does not require the knowledge of every scenario $\theta \in \Theta$ nor the scenario distribution $p$. Similar to empirical risk minimization (ERM) [13], a solution to (1) is obtained by solving the sample average approximation:

$$(\pi_{meta}, \phi) \in \arg\min \frac{1}{|\widehat{\Theta}|} \sum_{\theta \in \widehat{\Theta}} U_\theta(\phi(\pi, \theta)), \tag{2}$$

where $\widehat{\Theta} \subset \Theta$ is a finite collection of scenarios i.i.d. sampled from $p \in \Delta(\Theta)$. The term "agnostic" emphasizes that the exact scenario distribution $p$ is usually unknown in security practice and often replaced by an empirical distribution provided by security datasets, such as the data from MITRE ATT&CK [14] considered in the experiment section. Using the ERM language, it is expected that the empirical risk minimizer in (2) approximates the population risk minimizer in (1). When dealing with another scenario $\theta$ unseen in the sample set $\widehat{\Theta}$, the adapted policy $\phi(\pi_{meta}, \theta)$ achieves satisfying generalization.

Note that the pair $(\pi_{meta}, \phi)$ in (2) aims to minimize the sample average. However, a distribution shift between the empirical one and the true one in (1) can reduce the

adaptation/generalization ability, as the resulting meta policy may overfit to popular scenarios and perform poorly on rare ones [15], [16]. To address this issue, one can also design a distribution as presented in (3), leading to a robust ZTD that minimizes the worst-possible loss across all scenarios.

**Definition 2** (Scenario-Robust Zero-Trust Defense (SR-ZTD)). *A pair $\langle \pi_{meta}, \phi \rangle$ is scenario-robust if it solves for*

$$\min_{\pi, \phi} \sup_{p \in \Delta(\Theta)} \mathbb{E}_{\theta \sim p}[U_\theta(\phi(\pi, \theta))]. \tag{3}$$

The empirical approximation to (3) using $\widehat{\Theta}$ is given by $\min_{\pi, \phi} \max_{p \in \Delta(\widehat{\Theta})} \mathbb{E}_{\theta \sim p}[U_\theta(\phi(\pi, \theta))]$, which is equivalent to finding the optimum for the worst-possible case: $\min_{\pi, \phi} \max_{\theta \in \widehat{\Theta}} U_\theta(\phi(\pi, \theta))$, since $\Delta(\widehat{\Theta})$ is a probability simplex in a finite-dimensional space, and the support of the worst-case distribution contains some extreme points.

### A. Gradient-based Adaptation and Explainable Meta Policy

To elaborate on the adaptation mapping in SA-ZTD, we consider the following minimization problem in comparison to (1): $\min_\pi \mathbb{E}_{\theta \sim p}[U_\theta(\pi)]$. The corresponding minimizer is denoted by $\pi_{avg}$. When dealing with new tasks that are distant from the majority of training scenarios, $\pi_{avg}$ does not generalize well, as observed in Tables I and II. In contrast, the adaptation $\phi$ in SA-ZTD improves the generalization by updating $\pi_{meta}$ based on a few interactions in the new scenario, leading to a data-driven adaptation detailed below.

Since the function class $\{\phi | \phi : \Pi \times \Theta \to \Pi\}$ is infinite-dimensional, directly seeking an adaptation mapping through (1) [or (2)] is intractable. One remedy is to restrict the focus to the parameterization class where the mapping is parameterized by $\gamma \in \mathbb{R}^n$, $n \in \mathbb{Z}_+$. For example, $\phi_\gamma$ can be parameterized by recurrent neural networks, where $\gamma$ is the model weights, and the optimal adaptation is determined by training algorithms [17]. Another well-accepted parameterization is the gradient-based adaptation: $\phi_\gamma(\pi, \theta) := \pi - \gamma \nabla U_\theta(\pi)$, and $\gamma$ is the gradient step size to be optimized [18].

We consider the gradient-based adaptation due to its mathematical clarity and computational efficiency [19]. The adaptation step size $\gamma$ is assumed constant to simplify the exposition. The objective in SA-ZTD under gradient adaptation is $\min_\pi \mathbb{E}_{\theta \sim p}[U_\theta(\pi - \gamma \nabla U_\theta(\pi))]$ (fixing the adaptation), which implies that after obtaining the meta policy obtained through the meta training phase [i.e., solving (2)], one-step gradient descent update $\pi_{meta} - \gamma \nabla U_\theta(\pi_{meta})$ suffices for the new scenario when $\pi_{meta}$ is implemented. The gradient $\nabla U_\theta(\pi_{meta})$ in the deployment phase is estimated using a finite-different approximation scheme: simultaneous perturbation stochastic approximation (SPSA) [20], briefly reviewed in Algorithm 1. Such an estimate only requires finite horizon Monte Carlo (MC) simulations within the POMDP, leading to a lightweight adaptation without learning from scratch.

*a) Meta Threshold Policy:* Another reason to focus on gradient-based adaptation is that the resulting meta policy and adapted policy take threshold forms, leading to a switching

control. The active defense is implemented when the trust score (TS) is below the threshold.

To be specific, we first note that for each scenario $\theta$, the optimal policy, under some mild assumptions, is a stationary threshold policy [12, Theorem 1]. Hence, it suffices to find the minimizer $\pi_\theta^*$ to $U_\theta(\pi)$ within the set of threshold policies $\Pi$, where each element is parameterized by a threshold $\tau \in [0, 1]$: $\pi_\theta(TS) = \mathbb{1}_{\{\tau < TS \leq 1\}}$. Therefore, searching for the optimal policy is equivalent to finding the optimal threshold $\tau^*$. With a slight abuse of notations, we consider the cost $U_\theta(\pi)$ also a function of the threshold $\tau$, i.e., $U_\theta(\pi) = U_\theta(\tau)$. $U_\theta(\pi)$ and $U_\theta(\tau)$ are used interchangeably when no confusion arises.

To preserve the threshold form, we replace the one-step gradient adaptation with a projected gradient update: $\text{Proj}_{[0,1]}\{\tau - \gamma\nabla U_\theta(\tau)]\}$ so that the updated $\tau$ remains within $[0, 1]$, and the adapted policy is still a valid threshold policy. Hence, the objective in SA-ZTD (1) is now given by

$$\min_{\tau \in [0,1]} \mathbb{E}_{\theta \sim p}[U_\theta(\text{Proj}_{[0,1]}\{\tau - \gamma\nabla U_\theta(\tau)\})] \quad (4)$$

As a result, the meta policy takes the threshold form that is explainable to human operators, increasing the accessibility and transparency of learning-based ZTD. However, the price to pay is that (4) is a nonsmooth optimization problem, more involved than the original (1). To address this nonsmoothness, we propose an approximate stochastic gradient descent (SGD).

*b) Meta Learning Algorithm:* To solve (2) [and (3)], consider optimizing the empirical loss through SGD, i.e., $\tau_{meta}^{t+1} = \tau_{meta}^t - \frac{\alpha^t}{|\widehat{\Theta}_t|}\sum_{\theta \in \widehat{\Theta}_t}\nabla U_\theta(\phi(\tau_{meta}^t, \theta))\nabla_\tau\phi(\tau_{meta}^t, \theta)$, where $\alpha^t$ denotes the SGD step size, and $\widehat{\Theta}_t$ is a batch of scenarios sampled from $\widehat{\Theta}$ at $t$-th iteration. In addition to the non-smoothness issue, computing the gradient $\nabla_\tau\phi(\tau_{meta}^t, \theta)$ involves a Hessian matrix $\nabla^2 U_\theta$, and the Hessian-vector product computation is costly. To resolve these issues, we ignore term $\nabla_\tau\phi(\tau_{meta}^t, \theta)$ as commonly practiced in meta-learning applications [19], and it does not significantly affect the meta-learning performance as observed in [21].

Finally, it is numerically intractable to compute the exact policy gradient $\nabla U_\theta(\tau)$, and hence, we resort to the SPSA method (reviewed in Algorithm 1). Denote by $\widehat{\nabla} U_\theta(\tau)$ the SPSA gradient estimate under the policy $\pi$, then the one-step SGD update can be rewritten as

$$\tau_{meta}^{t+1} = \tau_{meta}^t - \frac{\alpha^t}{|\widehat{\Theta}_t|}\sum_{\theta \in \widehat{\Theta}_t}\widehat{\nabla} U_\theta(\tau_\theta^t), \quad (5)$$

$$\tau_\theta^t = \text{Proj}_{[0,1]}\{\tau_{meta}^t - \gamma\widehat{\nabla} U_\theta(\tau_{meta}^t)\}. \quad (6)$$

A summary of the above SPSA-based first-order meta-learning algorithm (SPSA-FOML) for SA-ZTD (1) is in Algorithm 1.

As for the minimax problem in SR-ZTD (3), we consider a stochastic gradient descent ascent (SGDA) algorithm, widely employed in adversarial meta learning [16], [22]. In SGDA, in addition to performing gradient descent on $\tau_{meta}^t$, a gradient ascent is applied to the variable $p^t \in \Delta(\widehat{\Theta})$: $p^{t+1}(\theta) = p^t(\theta) + \beta^t U_\theta(\phi(\tau_{meta}^t, \theta))$, for $\theta \in \widehat{\Theta}$, where $\beta^t$

is the ascent step size. Similar to the SGD update in SA-ZTD, the gradient adaptation $\phi(\tau_{meta}^t, \theta)$ is performed using SPSA estimate in (6). Meanwhile, the value function $U_\theta(\pi)$ is estimated through MC simulation by averaging the cumulative rewards of multiple finite-horizon MC rollouts [23]. Denote by $\widehat{U}_\theta$ the MC estimate, and for all $\theta \in \widehat{\Theta}^t$, the stochastic gradient ascent (SGA) update is given by

$$p^{t+1}(\theta) = p^t(\theta) + \beta^t\widehat{U}_\theta(\tau_\theta^t). \quad (7)$$

Note that after the SGA update, $p^{t+1}$ shall be projected to $\Delta(\widehat{\Theta})$. As one can see from Algorithm 1, SA-ZTD and SR-ZTD share the same meta-policy update using SGD [the minimization part in (3)], and the only difference is that SR-ZTD samples scenarios using the distribution updated by SGA.

*c) Complexity Analysis:* We briefly discuss the algorithmic complexity using results in [21] and [22]. Note that the bias of the SPSA gradient estimate is of $O[(\eta^t)^2]$ [20] (assuming summable). Hence, the summation $\sum_t \mathbb{E}[\|\nabla U_\theta(\tau_{meta}^t)\|^2]$ is controllable [22, Lemma E.5], which implies that SPSA-FOML converges to the $\epsilon$-first-order stationary point (line 14 in Algorithm 1) in SA-ZTD [21, Theorem 5.15] and SR-ZTD [22, Theorem E.2] within $O(\epsilon^{-2})$ iterations.

---

**Algorithm 1** SPSA-FOML for SA(SR)-ZTD

---

**Initialization**: threshold $\tau_{meta}^0 \in [0, 1]$, adaptation step size $\gamma$, scenario samples $\widehat{\Theta}$, scenario distribution $p^0 = \text{Uniform}(\widehat{\Theta})$, SPSA perturbation sequence $\{\eta^t\}$, gradient descent step sizes $\{\alpha^t\}$, ascent step sizes $\{\beta^t\}$, error tolerance $\epsilon$;

2: **for** $t = 0, 1, \ldots,$ **do**
    Sample a batch of scenarios $\widehat{\Theta}^t$ according to $p^t$;
4:    **for** every scenario $\theta \in \widehat{\Theta}^t$ **do**
        $\widehat{\nabla} U_\theta(\tau_{meta}^t) \leftarrow \text{SPSA}(\tau_{meta}^t, \eta^t, \theta)$;
6:        $\tau_\theta^t \leftarrow$ gradient-adaptation in (6);
        $\widehat{\nabla} U_\theta(\tau_\theta^t) \leftarrow \text{SPSA}(\tau_\theta^t, \eta^t, \theta)$;
8:        **if** SR-ZTD is activated **then**
            $\widehat{U}_\theta(\tau_\theta^t) \leftarrow$ MC simulation;
10:            $p^{t+1}(\theta) \leftarrow$ one-step SGA in (7)
        **else**
12:            $p^{t+1} \leftarrow \text{Uniform}(\widehat{\Theta})$
    $\tau_{meta}^{t+1} \leftarrow$ one-step SGD in (5)
14:    **if** $|\widehat{\nabla} U_\theta(\tau_{meta}^t)| \leq \epsilon$ for all $\theta$ **then** break
    **return** $\tau_{meta}^t$
16: **function** $\text{SPSA}(\tau, \eta, \theta)$         ▷ subroutine for SPSA
    $d \leftarrow \text{Bernoulli}(\{1, -1\}, 0.5)$;
18:    $\tau_+ \leftarrow \text{Proj}_{[0,1]}(\tau + \eta d)$; $\tau_- \leftarrow \text{Proj}_{[0,1]}(\tau - \eta d)$;
    $\widehat{U}_\theta(\tau_\pm) \leftarrow$ MC simulation using $\tau_\pm$;
20:    $\widehat{\nabla} U_\theta(\tau) \leftarrow \frac{\widehat{U}_\theta(\tau_+) - \widehat{U}_\theta(\tau_-)}{2\eta d}$
    **return** $\widehat{\nabla} U_\theta(\tau)$

---

## IV. EXPERIMENTAL RESULT

We present experimental results to show that the SA-ZTD policy is explainable based on scenario distributions and can adapt to the underlying scenario quickly. For illustrative purposes, we discuss two scenarios: varying system vulnerability $\theta = (p_u^n)$ and attacker's capability/stealthiness $\theta = (p_a^n)$.

**Experiment Setup**: The experiment parameters follow the suggestions in [20]: $\eta^t = 0.4/t^{0.2}$, $\alpha^t = \beta^t = 0.017/(t + 50)^{0.602}$, $\gamma = 0.005$, $\epsilon = 10^{-3}$. In the MC simulations, we

have $\rho = 0.86$, horizon length $T = 100$, $|\widehat{\Theta}| = 1000$ sampled from distributions to be defined later, and the batch size $|\widehat{\Theta}^t| = 10$. Each evaluation reports the mean and standard deviation under 50 repeated runs using different random seeds.

### A. System Vulnerability

We consider the following baseline configuration and train the model for different system vulnerability distributions: $p_a^d = 0.2, p_u^d = 0.1, p_a^n = 0.8, q_a = 0.9, q_u = 0.1, C(0,0) = 10, C(0,1) = 15, C(1,0) = 3, C(1,1) = 0$. To obtain the threshold policy [12], we consider the system vulnerability satisfies $0 \leq p_u^n \leq \min\{p_a^n, p_a^n - p_a^d + p_u^d\}$, i.e., $p_u^n \in [0, 0.7]$. Each $p_u^n$ value represents one scenario $\theta = (p_u^n)$ in our model. The larger $p_u^n$, the more vulnerable the system is.
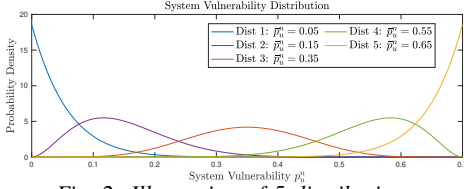


*Fig. 2: Illustration of 5 distributions.*

*1) Explainability:* We consider a generalized Beta distribution family $p_i(\theta, \alpha_i, \beta_i)$ with support $\theta = p_u^n \in [0, 0.7]$ to represent the vulnerability distributions. The $\alpha_i$ and $\beta_i$ values are chosen to generate means $\bar{p}_{u,i}^n = \mathbb{E}_{\theta \sim p_i}[\theta = p_u^n]$, whose values are shown in Figure 2.

In Figure 3, the black line illustrates the optimal policies under each fixed scenario $\theta = p_u^n$. We train the meta policy $\tau_{meta}$ under the 5 distributions and obtain the results, as in Figure 3. The dash lines are the single scenario optimal policies at $p_{u,min}^n$ and $p_{u,max}^n$. Experimental results agree with the explanation that if the system on average is more likely to be resistant to attacks (closer to $p_{u,min}^n$), $\tau_{meta}$ is relatively lower and closer to $\tau_{min}$. On the contrary, if the system on average is more vulnerable (closer to $p_{u,max}^n$), the meta policy focuses more on the vulnerable side and increases the minimum acceptable trust score.

*2) Adaptability:* In Figure 4, we look into the updated policy after a one-step gradient (also called one-shot) adaptation. The figure demonstrates that the adapted policies can capture the relationship between the system vulnerability and defense policy. The $\tau_{\theta,adapt}$ is lower when the actual scenario has lower $\theta = p_u^n$, while $\tau_{\theta,adapt}$ is increased if $\theta = p_u^n$ increases. It indicates that the meta policy can successfully adapt to each scenario. We also observe that when the training distribution average $\bar{p}_{u,i}^n$ is small (e.g., Dist. 1), the adaptation works better for the scenarios with small $p_u^n$. Similarly, when $\bar{p}_{u,i}^n$ is large (e.g., Dist. 5), the adaptation works better around high-vulnerability tasks.

To evaluate the defense performance, we randomly sample a finite collection of scenarios $\widehat{\Theta}_{test}$ i.i.d. from the underlying distribution (different from training sample $\widehat{\Theta}$). For each selected scenario, we use the one-shot adapted policy to evaluate our meta policy. We compare the average cost using the adapted meta policy $\overline{U}(\tau_{meta}) = 1/|\widehat{\Theta}_{test}| \sum_{\theta} U_\theta(\phi(\tau_{meta}, \theta))$ with the average cost always

*TABLE I: Adaptability with different system vulnerability.*

| Distribution | Average $\bar{p}_u^n$ | $\overline{U}(\pi_{meta})$ | $\overline{U}(\pi_{avg})$ |
|---|---|---|---|
| Dist. 1 | 0.05 | $18.98 \pm 1.01$ | $19.07 \pm 1.45$ |
| Dist. 2 | 0.15 | $22.69 \pm 1.50$ | $23.63 \pm 1.68$ |
| Dist. 3 | 0.35 | $29.10 \pm 1.54$ | $29.45 \pm 1.77$ |
| Dist. 4 | 0.55 | $30.63 \pm 1.70$ | $31.51 \pm 1.90$ |
| Dist. 5 | 0.65 | $31.25 \pm 1.23$ | $32.01 \pm 1.41$ |

*TABLE II: Adaptability with different attacker's stealthiness.*

| Distribution | Average $\bar{p}_a^n$ | $\overline{U}(\pi_{meta})$ | $\overline{U}(\pi_{avg})$ |
|---|---|---|---|
| Dist. 1 | 0.65 | $25.46 \pm 0.41$ | $26.01 \pm 0.59$ |
| Dist. 2 | 0.7 | $27.89 \pm 0.54$ | $28.40 \pm 0.77$ |
| Dist. 3 | 0.8 | $29.08 \pm 0.53$ | $31.10 \pm 0.55$ |
| Dist. 4 | 0.9 | $32.98 \pm 0.79$ | $33.52 \pm 1.03$ |
| Dist. 5 | 0.95 | $34.15 \pm 0.74$ | $34.23 \pm 0.57$ |

using the optimal policy at distribution average $\overline{U}(\tau_{avg}) = 1/|\widehat{\Theta}_{test}| \sum_{\theta} U_\theta(\tau_{avg})$. It should be noted that there exist testing scenarios that the defender had not encountered during the training process.

Table I compares the average cost between the adapted meta policy and distribution average policy. We observe that our meta policy outperforms the distribution-average policy under all 5 distributions as both the average cost and standard deviation are lower. It indicates that the meta policy has the adaptation ability to the sampled scenarios and provides a lower average cost compared to always using the distribution average policy. SA-ZTD can provide a tailored defense policy against different system configurations.

### B. Attacker's Stealthiness

To illustrate the results, we keep the baseline settings and let $p_u^n = 0.5$. The attacker's stealthiness is captured by $p_a^n$ and we consider $\max\{p_u^n, p_u^n - p_u^d + p_a^d\} \leq p_a^n \leq 1$, i.e., $p_a^n \in [0.6, 1]$, to keep the threshold policy [12]. Each $p_a^n$ value represents one scenario $\theta = (p_a^n)$ in the model.

*1) Explainability:* To evaluate the meta policy under different attacker distributions, again, we consider 5 generalized Beta distributions $p_i'(\theta, \alpha_i, \beta_i)$ on support $\theta = p_a^n \in [0.6, 1]$ with mean values $\bar{p}_{a,i}^n \in \{0.65, 0.7, 0.8, 0.9, 0.95\}$, respectively. The meta policy is summarized in Figure 5. The single scenario optimal policy indicates that the system must increase the trust threshold $\tau$ if the attacker is more capable. As shown in the figure, the adapted policies can capture the relationship between the attacker's stealthiness and defense policy. When $\bar{p}_a^n$ is closer to $p_{a,min}^n$, the meta policy is closer to $\tau_{min}$. When $\bar{p}_a^n$ is closer to $p_{a,max}^n$, the meta policy is also closer to $\tau_{max}$. The observations concur with the common explanation that if the attacker is more likely to be stealthy, the defense system needs to guard up and increase the meta-trust threshold. Conversely, if the attacker is less capable, the zero-trust engine could be more tolerant of the agent with a lower trust score.

*2) Adaptability:* We investigate the adaptation ability of our policy to different attackers' stealthiness levels in Figure 6. After adaptation, the updated policy successfully adapts to each underlying scenario as it generates a lower trust threshold with small $p_a^n$ and provides a higher threshold with large $p_a^n$.

Finally, we randomly sample a finite set of scenarios $\widehat{\Theta}_{test}$ i.i.d. from the underlying distribution and compare $\overline{U}(\tau_{meta})$ with $\overline{U}(\tau_{avg})$. The results in Table II demonstrate that the meta policy generates a lower cost with less variance compared
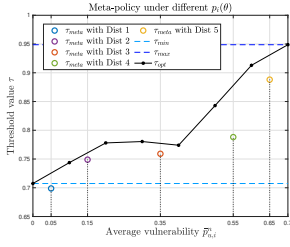
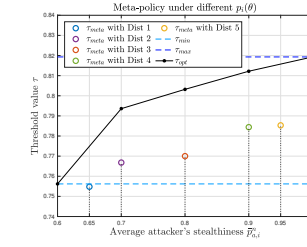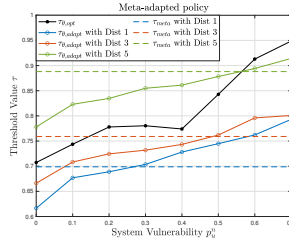Fig. 3: Meta policy under different system vulnerability distributions. Fig. 4: Adapted policy $\tau_{adapt}$ after one-step gradient update. Fig. 5: Meta policy under different attacker stealthiness distributions. Fig. 6: Adapted policy $\tau_{adapt}$ after one-step gradient update.

to the distribution average cost under all 5 distributions. Our method can quickly adapt to the specific scenario and provide a utility-wise better defense policy.
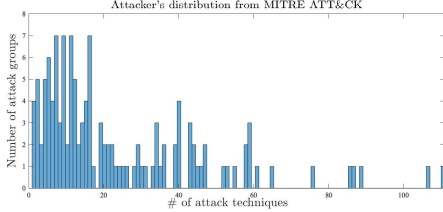
### C. Real Attack Data



Fig. 7: Attacker's distribution histogram from MITRE ATT&CK [14].

We finally evaluate SA-ZTD on real-world attack distribution. We consider the scenarios $\theta = (p_a^n)$ and collect real-world attack group data from MITRE ATT&CK [14]. The histogram of the attacker's distribution $\Delta(\Theta)$ is illustrated in Figure 7, where the x-axis is the number of attack techniques and the y-axis is the number of attack groups. We use this empirical distribution to approximate the distribution of $p_a^n$ and compute the meta policy.

TABLE III: System performance with different policies.

| Distribution | $\overline{U}(\pi_{avg})$ | $\overline{U}(\pi_{meta})$ | $\overline{U}(\pi_{robust})$ |
|---|---|---|---|
| Empirical Dist. | $28.89 \pm 1.11$ | $\mathbf{27.89} \pm 0.66$ | $32.11 \pm 0.95$ |
| Worst-case Dist. | $36.18 \pm 0.86$ | $33.75 \pm 1.51$ | $\mathbf{32.48} \pm 1.39$ |

We randomly sample a set of scenarios from empirical and worst-case distributions and compare the average costs with different policies. Table III shows that the meta policy SA-ZTD outperforms the other two cases when we use empirical distribution while the robust policy SR-ZTD has a lower cost when we consider the worst-case distribution.

### V. CONCLUSION

We have formulated a scenario-agnostic zero-trust defense (SA-ZTD) and obtained an explainable and adaptable trust-threshold policy that activates defense based on the agent's trust score. We have developed an algorithm using first-order meta-learning to learn the SA-ZTD policy and extend it to robust defense in response to real-world data. Experiments have demonstrated the explainability and adaptability of the proposed model.

### REFERENCES

[1] Q. Zhu and Z. Xu, *Cross-Layer Design for Secure and Resilient Cyber-Physical Systems*. Springer, 2020.

[2] Q. Zhu, S. Rass, B. Dieber, and V. M. Vilches, "Cybersecurity in robotics: Challenges, quantitative modeling, and practice," *arXiv preprint arXiv:2103.05789*, 2021.

[3] S. Rose, O. Borchert, S. Mitchell, and S. Connelly, "Zero trust architecture," tech. rep., National Institute of Standards and Technology, 2020.

[4] T. M. Hospedales, A. Antoniou, P. Micaelli, and A. J. Storkey, "Meta-Learning in Neural Networks: A Survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, no. 99, pp. 1–1, 2021.

[5] Q. Zhu and Z. Xu, "Introduction to partially observed mdps," in *Cross-Layer Design for Secure and Resilient Cyber-Physical Systems*, pp. 139–145, Springer, 2020.

[6] S. Mehraj and M. T. Banday, "Establishing a zero trust strategy in cloud computing environment," in *2020 International Conference on Computer Communication and Informatics (ICCCI)*, pp. 1–6, IEEE, 2020.

[7] K. Ramezanpour and J. Jagannath, "Intelligent zero trust architecture for 5g/6g networks: Principles, challenges, and the role of machine learning in the context of o-ran," *Computer Networks*, p. 109358, 2022.

[8] S. Dhar and I. Bose, "Securing iot devices using zero trust and blockchain," *Journal of Organizational Computing and Electronic Commerce*, vol. 31, no. 1, pp. 18–34, 2021.

[9] R. Heartfield and G. Loukas, "A taxonomy of attacks and a survey of defence mechanisms for semantic social engineering attacks," *ACM Computing Surveys (CSUR)*, vol. 48, no. 3, pp. 1–39, 2015.

[10] L. Ablon and A. Bogart, *Zero days, thousands of nights: The life and times of zero-day vulnerabilities and their exploits*. Rand Corporation, 2017.

[11] T. Li, Y. Zhao, and Q. Zhu, "The role of information structures in game-theoretic multi-agent learning," *Annual Reviews in Control*, vol. 53, pp. 296–314, 2022.

[12] Y. Ge and Q. Zhu, "Trust Threshold Policy for Explainable and Adaptive Zero-Trust Defense in Enterprise Networks," *2022 IEEE Conference on Communications and Network Security (CNS)*, pp. 359–364, 2022.

[13] V. Vapnik, *The nature of statistical learning theory*. Springer science & business media, 1999.

[14] B. E. Strom, A. Applebaum, D. P. Miller, K. C. Nickels, A. G. Pennington, and C. B. Thomas, "Mitre att&ck: Design and philosophy," in *Technical report*, The MITRE Corporation, 2018.

[15] B. Mehta, T. Deleu, S. C. Raparthy, C. J. Pal, and L. Paull, "Curriculum in Gradient-Based Meta-Reinforcement Learning," *arXiv*, 2020.

[16] L. Collins, A. Mokhtari, and S. Shakkottai, "Task-robust model-agnostic meta-learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 18860–18871, 2020.

[17] S. Y. Hochreiter, "Learning to Learn Using Gradient Descent," *Lecture Notes in Computer Science*, pp. 87–94, 2001.

[18] Z. Li, F. Zhou, F. Chen, and H. Li, "Meta-SGD: Learning to Learn Quickly for Few-Shot Learning," *arXiv*, 2017.

[19] C. Finn, P. Abbeel, and S. Levine, "Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks," *arXiv*, 2017.

[20] J. C. Spall, *Introduction to stochastic search and optimization: estimation, simulation, and control*. John Wiley & Sons, 2005.

[21] A. Fallah, A. Mokhtari, and A. Ozdaglar, "On the convergence theory of gradient-based model-agnostic meta-learning algorithms," in *International Conference on Artificial Intelligence and Statistics*, pp. 1082–1092, PMLR, 2020.

[22] T. Li, H. Lei, and Q. Zhu, "Sampling Attacks on Meta Reinforcement Learning: A Minimax Formulation and Complexity Analysis," *arXiv*, 2022.

[23] D. P. Bertsekas and J. N. Tsitsiklis, "Neuro-dynamic programming," Optimization and neural computation series, 1996.