

Divergence-Minimizing Attack Against Challenge-Response Authentication with IRSs

Laura Crosara, Anna V. Guglielmi, Nicola Laurenti, and Stefano Tomasin

Dept. of Information Engineering (DEI), University of Padova, Italy

email: {laura.crosara.l@phd., annavaleria.guglielmi@, stefano.tomasin@, nicola.laurenti@}unipd.it

Abstract—We propose a new attack against challenge response physical layer authentication (CR-PLA) with intelligent reflecting surfaces (IRSs). Drawing from prior work, we establish bounds on performance metrics, such as probabilities of false alarm and missed detection, using Kullback-Leibler (KL) divergence. Leveraging prior results in [1], we extend the analysis to the CR-PLA scenario with IRS. We derive the optimal attack strategy to minimize the divergence between authentic and forged signals when the attacker has either partial or no knowledge of the legitimate cascade channel. We evaluate the attack performance under different conditions, by varying the correlation between Eve’s observation and the legitimate cascade channel, or the SNR at the legitimate receiver.

Index Terms—Authentication, Challenge-response, Intelligent Reflecting Surfaces.

I. INTRODUCTION

Source authentication is the problem of establishing if a received message truly comes from the declared sender or has been forged by an impersonating attacker. Risks in accepting unauthenticated messages go from denial of service to privacy to the loss of control of devices, e.g., in Internet of Things (IoT) contexts.

Several authentication mechanisms have been explored, beyond those operating at the application level and using cryptographic approaches. Here we focus on physical layer authentication (PLA), which leverages the propagation characteristics of the physical channel as signatures of the transmitting devices or the communication links. The basic approach, introduced by Simmons [2] includes two phases, the identification acquisition, and the identification verification. In the former, the receiver Bob (verifier) estimates the channel using signals transmitted by Alice (the authentic source) that are authenticated at higher layers (e.g., by cryptographic mechanisms). In the latter, whenever Bob receives a new message, he also estimates the channel over which the signal traveled and compares this estimate with that obtained in the first phase. If the two estimates are consistent (note that they are still affected by noise), then the received message is deemed authentic, otherwise it is considered fake. PLA has been studied for several technologies, including orthogonal frequency

This work has been funded in part by the European Commission through the Horizon Europe/JU SNS project ROBUST-6G (Grant Agreement no. 101139068).

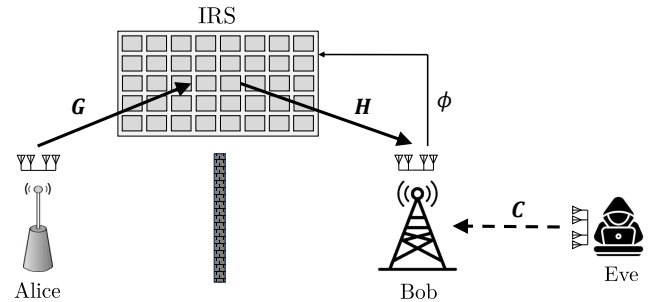


Fig. 1. Communication scenario.

division multiplexing (OFDM) and multiple-input multiple-output (MIMO) [3], [4], underwater acoustic communications [5], and using several techniques for the test, from Neyman-Pearson tests [6] to machine learning approaches [7]. For an overview on PLA see also [8], [9].

Recently, a further evolution of PLA has been introduced by exploiting the controllable nature of wireless channels provided by new communication technologies. In particular, intelligent reflecting surfaces (IRSs) are controllable devices that can change the propagation of wireless signals by changing the phase shift introduced by their elements. When the IRS is under Bob’s control, he can set a random configuration of the IRS which remains secret to the attacker, and verify that the estimated channel on a received message corresponds to the set configuration, [10]. This approach provides a *challenge response PLA (CR-PLA)* mechanism and can be applied also when other *controllable channels* are available (e.g., Bob is a drone that can change its position, [11]).

In this paper, we investigate novel attacks to be performed when the CR-PLA uses an IRS to perform the challenge-response approach. In particular, we leverage the results of [1, Th. 2] that has established bounds on the performance for a conventional (non-interactive) PLA mechanism, in terms of probabilities of false alarm (FA) and missed detection (MD). The bounds exploit the Kullback-Leibler (KL) divergence of observed channels at Alice, Bob, and Eve, and it turns out that when the observations in the legitimate case are jointly Gaussian distributed, the optimal attack strategy is also Gaussian distributed. Here we derive the bounds for the considered CR-PLA scenario with IRS and, under the assumption of a large number of IRS elements, we ensure that the assumption

in [1, Th. 2] are satisfied and derive the optimal attack by Eve. We assess the performance of the obtained attack by considering different correlations between Eve's observation and the legitimate cascade channel, and different signal-to-noise ratio (SNR) values for the legitimate channel.

The rest of the paper is organized as follows. Section II presents the system model. Section III describes the CR-PLA mechanism. Then, the KL divergence minimizing attack is derived in Section IV, considering no channel knowledge at Eve in Section IV-A, and partial channel knowledge in Section IV-B. Then, numerical results are discussed in Section V and, finally, conclusions are drawn in Section VI.¹

II. SYSTEM MODEL

We consider the scenario depicted in Fig. 1, where a legitimate receiver Bob authenticates messages from a legitimate transmitter Alice. We account for the presence of an attacker Eve, who aims to impersonate Alice by forging messages and transmitting them to Bob. We assume that Alice, Bob, and Eve are equipped with a uniform linear array (ULA) of N_A , N_B , and N_E antennas, respectively.

We further assume that the communication between Alice and Bob is supported by an IRS with N reflecting elements, each acting as a receive and transmit antenna. The n -th element, $n = 1, 2, \dots, N$, of the IRS introduces a phase shift $\phi_n = e^{j\theta_n}$ on the equivalent baseband signal and has unitary gain. We define vector $\phi = [\phi_1, \phi_2, \dots, \phi_N]^T$ as the *configuration* that collects the phase shifts introduced by the IRS on each element, and the $N \times N$ diagonal matrix

$$\Phi = \text{diag}\{\phi\} = \text{diag}\{\phi_1, \phi_2, \dots, \phi_N\}. \quad (1)$$

Bob controls the IRS and sets ϕ using a secure dedicated channel (between Bob and the IRS) not accessible to Eve. We assume that communication between Alice and Bob only happens through the IRS without any additional direct link (for instance, because a direct link is not available). We denote the baseband equivalent matrix for the channel from Alice to the IRS as $\mathbf{G} \in \mathbb{C}^{N \times N_A}$, the channel from the IRS to Bob as $\mathbf{H} \in \mathbb{C}^{N_B \times N}$. Thus, the resulting Alice-IRS-Bob cascade channel is

$$\mathbf{Q}_{\text{AIB}} = \mathbf{H}\Phi\mathbf{G}. \quad (2)$$

We assume that Eve can directly transmit to Bob through the channel \mathbf{C} , whose channel matrix is known by Eve.

Note that all channels are time-invariant and reciprocal, while the IRS configuration that reflects onto matrix Φ is under Bob's control and can be changed over time, making the cascade channels controllable.

¹Notation: \dagger denotes the Moore-Penrose inverse, H represents complex conjugate and transpose, $|\mathbf{A}|$ stands for the determinant of \mathbf{A} , and \log is the natural logarithm. If $\mathbf{a} \in \mathbb{C}^n$ and $\mathbf{b} \in \mathbb{C}^m$ are random vectors, $\mathbf{K}_{ab} = \mathbb{E}[\mathbf{a}\mathbf{b}^H]$ is their $n \times m$ covariance matrix. The vectorization operator $[\text{vec}(\mathbf{X})]_i = [\mathbf{X}]_{\lfloor \frac{i}{RC} \rfloor, i \bmod RC}$, $i = 1, \dots, RC$, converts matrix \mathbf{X} of size $(R \times C)$ to a column vector of length RC by indexing the matrix by rows.

III. CHALLENGE-RESPONSE PLA WITH IRS

The CR-PLA mechanism introduced in [10] works as follows. In the *association* phase Alice transmits several authenticated pilot signals and Bob estimates the cascade channel \mathbf{Q}_{AIB} for various IRS configurations. We assume that this phase is long enough to let Bob obtain estimates $\hat{\mathbf{Q}}_{\text{AIB}}(\Phi')$ of the cascade channel for any IRS configuration Φ' .

In the *verification* phase, Bob sets a random configuration Φ of the IRS as follows. Bob modifies the IRS configuration around the communication-optimal² one $\bar{\theta}_n$, $n = 1, \dots, N$, obtaining a new configuration as

$$\phi'_n = e^{j(\bar{\theta}_n + \epsilon_n)}, \quad n = 1, \dots, N, \quad (3)$$

with $\epsilon_n \sim \mathcal{U}[-\gamma, \gamma]$, where γ denotes the maximum deviation from the communication-optimal IRS configuration. Thus, a larger γ introduces more randomness in the IRS configuration making it more difficult for Eve to build an attack that passes the authentication check. Moreover, ϵ_n are independent for each $n = 1, \dots, N$, and regenerated independently at each new verification phase. Whenever Bob receives a message, he estimates the cascade channel \mathbf{R} and checks if it is consistent with the expected channel $\mathbf{Q}_{\text{AIB}}(\Phi)$. Pilot signals are assumed to be known by all parties.

Under legitimate conditions, given the configuration ϕ , when Alice transmits, Bob obtains the channel estimate

$$\hat{\mathbf{Q}}_{\text{AIB}} = \mathbf{Q}_{\text{AIB}} + \mathbf{W}_B, \quad (4)$$

where \mathbf{W}_B denotes the estimation error at Bob, modeled as additive white Gaussian noise (AWGN) with zero mean, independent entries and power σ_B^2 per entry.

We assume that Bob has perfect knowledge of both channel matrices and can choose the IRS configuration Φ . Thus, Bob during the association phase obtains $\hat{\mathbf{Q}}_{\text{AIB}}(\Phi) = \mathbf{Q}_{\text{AIB}}(\Phi)$.

Bob's goal is to figure out whether the estimated channel \mathbf{R} is authentic ($\hat{\mathbf{Q}}_{\text{AIB}}$), or forged (\mathbf{V}), by using his knowledge of \mathbf{Q}_{AIB} . To detect the attack, Bob performs an authentication test, that, given \mathbf{Q}_{AIB} and \mathbf{R} , decides between the hypotheses:

$$\mathcal{H}_0 : \text{the message is from Alice}, \quad (5)$$

$$\mathcal{H}_1 : \text{the message is from Eve}. \quad (6)$$

The authentication procedure is summarized in block D of Fig. 2, which has \mathbf{r} as input and outputs the Boolean value \hat{b} . Correct verification is achieved when $\hat{b} = b$.

A. Attack Model

The purpose of Eve is to impersonate Alice, i.e., to pass the authentication check at the legitimate receiver. We consider that some side information is available to Eve, represented by the matrix \mathbf{Z} . We assume that \mathbf{C} is correlated with \mathbf{H} and/or \mathbf{G} , and Eve can estimate these channel matrices and compute (an estimate of) the communication-optimal IRS configuration $\Phi^* = \text{diag}\{e^{j\bar{\theta}_1}, \dots, e^{j\bar{\theta}_N}\}$. A simple model of the resulting

²For further insights on optimal IRS configuration refer to [12].

estimate at Eve of the Alice-IRS-Bob cascade channel with communication-optimal IRS configuration is

$$\mathbf{Z} = \mathbf{H}\Phi^*\mathbf{G} + \mathbf{W}_z, \quad (7)$$

where \mathbf{W}_z is an AWGN matrix modeling the resulting estimation error. In Section IV we will analyze the case where Eve's side information \mathbf{Z} and the legitimate cascade channel \mathbf{X} are either independent or correlated.

We assume that Eve knows: i) the actual realization of \mathbf{Z} , ii) the probability density function (pdf) of the channel matrices \mathbf{H} and \mathbf{G} , and of the IRS configuration vector ϕ , and iii) the pdf of noise at both receivers. We assume that Eve transmits precoded pilot signals and can induce any channel estimate (apart from the presence of noise) to Bob. Then, the channel estimated by Bob when under attack is

$$\mathbf{V} = \mathbf{V}_0 + \mathbf{W}_B, \quad (8)$$

where \mathbf{V}_0 is the channel selected by Eve for the attack. For the sake of generality, we consider that Eve adopts a probabilistic strategy, characterized by the conditional pdf $p_{V_0|\mathbf{Z}}$. Moreover, since Eve knows the statistics of \mathbf{W}_B , the attack strategy can be described by $p_{V|\mathbf{Z}}$.

Finally, let the channel estimated at Bob be

$$\mathbf{R} = \begin{cases} \hat{\mathbf{Q}}_{\text{AIB}} & \text{if Alice is transmitting } (b = 0), \\ \mathbf{V} & \text{if Eve is transmitting } (b = 1), \end{cases} \quad (9)$$

where b indicates the legitimate/attack state. Therefore, Eve aims to prevent Bob from distinguishing between the attack \mathbf{V} and the legitimate $\hat{\mathbf{Q}}_{\text{AIB}}$. For ease of notation, let $\mathbf{v} = \text{vec}(\mathbf{V})$ and $\mathbf{r} = \text{vec}(\mathbf{R})$, as also shown in Fig. 2.

From now on, to compact the notation, we let

$$\mathbf{x} = \text{vec}(\hat{\mathbf{Q}}_{\text{AIB}}), \quad (10)$$

$$\mathbf{y} = \text{vec}(\hat{\mathbf{Q}}_{\text{AIB}}) = \mathbf{x} + \boldsymbol{\omega}, \quad (11)$$

$$\mathbf{z} = \text{vec}(\mathbf{Z}), \quad (12)$$

where $\boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_\omega)$, $\mathbf{K}_\omega = \sigma_B^2 \mathbf{I}$. So, \mathbf{x} is the information available to Bob during the verification phase, \mathbf{y} is the legitimate channel estimation, and \mathbf{z} contains the information Eve can leverage to build the attack. Moreover, we refer to the abstract representation of the authentication scenario in Fig. 2.

B. Performance Metrics

An authentication system's performance is typically evaluated through type-I (FA) error probability p_{FA} , i.e., the probability that Bob rejects a legitimate message from Alice, and type-II (MD) error probability p_{MD} , i.e., the probability that Bob accepts a message from Eve as authentic. The analytical derivation of such probabilities, depending on the verification decision rule, is often impractical. Thus, establishing bounds on the achievable error probability region, i.e., the set of feasible points in the $(p_{\text{FA}}, p_{\text{MD}})$ plane, is useful to prove the effectiveness of practical schemes.

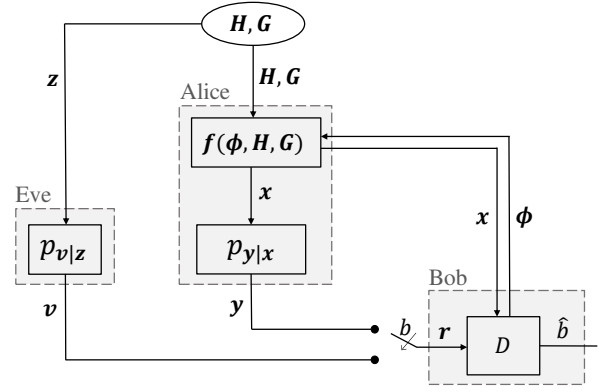


Fig. 2. Block scheme.

An outer bound is obtained by the KL divergence among the estimated channels at Alice, Bob, and Eve. From [6], as a consequence of the data processing inequality, we have

$$\mathbb{D}(p_{\hat{b}|\mathcal{H}_1} \| p_{\hat{b}|\mathcal{H}_0}) \leq \mathbb{D}(p_{\mathbf{r}|\mathcal{H}_1} \| p_{\mathbf{r}|\mathcal{H}_0}) = \mathbb{D}(p_{\mathbf{x}v} \| p_{\mathbf{x}y}), \quad (13)$$

where $\mathbb{D}(p \| q)$ is the KL divergence between pdfs p and q . It also holds the symmetric bound

$$\mathbb{D}(p_{\hat{b}|\mathcal{H}_0} \| p_{\hat{b}|\mathcal{H}_1}) \leq \mathbb{D}(p_{\mathbf{r}|\mathcal{H}_0} \| p_{\mathbf{r}|\mathcal{H}_1}) = \mathbb{D}(p_{\mathbf{x}y} \| p_{\mathbf{x}v}). \quad (14)$$

In (13) and (14) we have considered $p_{\mathbf{r}x}$ based on the assumption that, at the time of verification, the legitimate receiver knows \mathbf{x} , and the decision \hat{b} is made on both inputs. Therefore, observing that $p_{\hat{b}|\mathcal{H}_0}(1) = p_{\text{FA}}$, $p_{\hat{b}|\mathcal{H}_1}(0) = p_{\text{MD}}$, and defining the function

$$h(p_{\text{MD}}, p_{\text{FA}}) \triangleq p_{\text{MD}} \log \frac{p_{\text{MD}}}{1 - p_{\text{FA}}} + (1 - p_{\text{MD}}) \log \frac{1 - p_{\text{MD}}}{p_{\text{FA}}}, \quad (15)$$

(13) can be rewritten as

$$h(p_{\text{MD}}, p_{\text{FA}}) \leq \mathbb{D}(p_{\mathbf{x}v} \| p_{\mathbf{x}y}). \quad (16)$$

This limits the region of achievable $(p_{\text{FA}}, p_{\text{MD}})$ values, depending on $\mathbb{D}(p_{\mathbf{x}v} \| p_{\mathbf{x}y})$, i.e., the statistics of the observations at Bob and Eve, for any decision mechanism. The use of challenge-response authentication has an impact on \mathbf{x} (and \mathbf{y}) through the additional randomness introduced by the choice of the IRS configuration in the verification phase. Since this randomness is not captured by \mathbf{z} , the challenge-response mechanism increases the KL divergence at Bob's benefit. Thus, a higher value of γ in the model (3) for the IRS randomness will yield a larger KL divergence.

The aim of the attacker Eve is to narrow the achievable region, by making the value of $\mathbb{D}(p_{\mathbf{x}v} \| p_{\mathbf{x}y})$ as small as possible, operating on the attack strategy $p_{v|\mathbf{z}}$. The metric $\mathbb{D}(p_{\mathbf{x}v} \| p_{\mathbf{x}y})$ depends on the attack strategy as

$$\mathbb{D}(p_{\mathbf{x}v} \| p_{\mathbf{x}y}) = \mathbb{E} \left[\log \frac{\int p_{v|\mathbf{z}}(v|\mathbf{a}) p_{z|\mathbf{x}}(\mathbf{a}|\mathbf{x}) d\mathbf{a}}{p_{y|\mathbf{x}}(v|\mathbf{x})} \right]. \quad (17)$$

IV. DIVERGENCE MINIMIZING ATTACK

In this Section, we derive the attack strategy that minimizes the KL divergence, thus abstracting from the particular detection process, whose efficacy can be assessed beforehand. We consider two scenarios: 1) Eve has no information about the Alice-IRS-Bob channel, i.e., \mathbf{z} and \mathbf{x} are independent; 2) Eve has partial information about the Alice-IRS-Bob channel, thus \mathbf{z} is correlated with \mathbf{x} .

Now, assuming that the number of IRS elements N is large enough, since each element $x_{p,q}$ of the matrix \mathbf{Q}_{AIB} is given by $x_{p,q} = \sum_{i=1}^N \phi_i \mathbf{H}_{p,i} \mathbf{G}_{i,q}$, considering the model (3) for the IRS configuration and invoking the central limit theorem, we can model \mathbf{x} as a Gaussian distributed random vector with vector mean $\boldsymbol{\mu}_x$ and covariance matrix \mathbf{K}_x . So, \mathbf{x} and \mathbf{y} are jointly Gaussian random vectors with covariance matrix

$$\mathbf{K}_{\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}} = \begin{bmatrix} \mathbf{K}_x & \mathbf{K}_{xy} \\ \mathbf{K}_{yx} & \mathbf{K}_y \end{bmatrix} = \begin{bmatrix} \mathbf{K}_x & \mathbf{K}_x \\ \mathbf{K}_x & \mathbf{K}_x + \mathbf{K}_\omega \end{bmatrix}. \quad (18)$$

A. No Channel Knowledge

In the zero channel knowledge hypothesis, \mathbf{x} and \mathbf{v} are independent, thus, $p_{xv} = p_x p_v$. Under these conditions, $\mathbb{D}(p_{xy} \| p_{xv})$ becomes tractable, hence we consider two cases for the design of the attack: *a)* Eve aims at minimizing $\mathbb{D}(p_{xv} \| p_{xy})$ or *b)* Eve aims at minimizing $\mathbb{D}(p_{xy} \| p_{xv})$.

Considering task *a)*, the optimal joint pdf p_{xv} that minimizes $\mathbb{D}(p_{xv} \| p_{xy})$ is also Gaussian [1, Th. 2]. When the attack strategy is $p_v \sim \mathcal{N}(\boldsymbol{\mu}_v, \mathbf{K}_v)$, we have

$$\begin{aligned} \mathbb{D}(p_x p_v \| p_{xy}) &= \mathbb{E} \left[\log \frac{p_v(\mathbf{v})}{p_{y|x}(\mathbf{v}|\mathbf{x})} \right] = \mathbb{E} \left[\log \frac{p_v(\mathbf{v})}{p_\omega(\mathbf{v} - \mathbf{x})} \right] \\ &= \mathbb{E} \left[\log \frac{p_v(\mathbf{v})}{p_\omega(\mathbf{v})} \right] + \mathbb{E} \left[\log \frac{p_\omega(\mathbf{v})}{p_\omega(\mathbf{v} - \mathbf{x})} \right] \end{aligned} \quad (19)$$

The first term in (19) is the KL divergence $\mathbb{D}(p_v \| p_\omega)$, while the second term can be computed as

$$\begin{aligned} \mathbb{E} \left[\log \frac{p_\omega(\mathbf{v})}{p_\omega(\mathbf{v} - \mathbf{x})} \right] &= \frac{1}{2} \left[\text{tr}(\mathbf{K}_x \mathbf{K}_\omega^{-1}) \right. \\ &\quad \left. - \boldsymbol{\mu}_x^H \mathbf{K}_\omega^{-1} \boldsymbol{\mu}_v - \boldsymbol{\mu}_v^H \mathbf{K}_\omega^{-1} \boldsymbol{\mu}_x \right]. \end{aligned} \quad (20)$$

Replacing (20) in (19) we obtain

$$\begin{aligned} \mathbb{D}(p_x p_v \| p_{xy}) &= \frac{1}{2} \left[\log \frac{|\mathbf{K}_\omega|}{|\mathbf{K}_v|} + \text{tr}(\mathbf{K}_v \mathbf{K}_\omega^{-1}) - N_A N_B \right. \\ &\quad \left. + \text{tr}(\mathbf{K}_x \mathbf{K}_\omega^{-1}) + (\boldsymbol{\mu}_v - \boldsymbol{\mu}_x)^H \mathbf{K}_\omega^{-1} (\boldsymbol{\mu}_v - \boldsymbol{\mu}_x) \right], \end{aligned} \quad (21)$$

which is minimized for $\boldsymbol{\mu}_v = \boldsymbol{\mu}_x$, $\mathbf{K}_v = \mathbf{K}_\omega$. So, when Eve has no clue about the cascade channel Alice-IRS-Bob the optimal probabilistic attack to minimize $\mathbb{D}(p_{xv} \| p_{xy})$ is

$$p_v^* \sim \mathcal{N}(\boldsymbol{\mu}_x, \mathbf{K}_\omega) \quad (22)$$

where Eve mimics the estimation error distribution at Bob, centered at $\boldsymbol{\mu}_x$. When Eve strategy is p_v^* , (19) becomes

$$\mathbb{D}(p_x p_v^* \| p_{xy}) = \frac{1}{2} \text{tr}(\mathbf{K}_x \mathbf{K}_\omega^{-1}) = \frac{1}{2} \sigma_B^{-2} \text{tr}(\mathbf{K}_x). \quad (23)$$

Moreover, with this attack strategy, we also get $\mathbb{D}(p_x p_v^* \| p_{xy}) = \mathbb{D}(p_{xy} \| p_x p_v^*)$.

Considering now case *b)*, Eve aims to find the pdf p_v that minimizes $\mathbb{D}(p_{xy} \| p_{xv})$. For any joint pdf p_{xy} , even non Gaussian, we have $\mathbb{D}(p_{xy} \| p_{xv}) = \mathbb{D}(p_{xy} \| p_x p_v)$, and

$$\begin{aligned} \mathbb{D}(p_{xy} \| p_x p_v) &= \\ &= \mathbb{E} \left[\log \frac{p_{xy}(\mathbf{x}, \mathbf{y})}{p_x(\mathbf{x}) p_v(\mathbf{y})} \right] = \mathbb{E} \left[\log \frac{p_{xy}(\mathbf{x}, \mathbf{y})}{p_x(\mathbf{x}) p_y(\mathbf{y})} \frac{p_y(\mathbf{y})}{p_v(\mathbf{y})} \right] \quad (24) \\ &= \mathbb{I}(\mathbf{x}; \mathbf{y}) + \mathbb{D}(p_y \| p_v) \geq \mathbb{I}(\mathbf{x}; \mathbf{y}), \end{aligned}$$

where $\mathbb{I}(\mathbf{x}; \mathbf{y})$ denotes the mutual information between random variables \mathbf{x} and \mathbf{y} . Therefore, the optimal probabilistic attack strategy is

$$p_v^* = \arg \min_{p_v} \mathbb{D}(p_{xy} \| p_x p_v) = \arg \min_{p_v} \mathbb{D}(p_y \| p_v) = p_y \quad (25)$$

and the KL divergence in the jointly Gaussian case (18) is

$$\mathbb{D}(p_{xy} \| p_x p_v^*) = \mathbb{I}(\mathbf{x}; \mathbf{y}) = \frac{1}{2} \log \left(\frac{|\mathbf{K}_x + \mathbf{K}_\omega|}{|\mathbf{K}_\omega|} \right). \quad (26)$$

B. Partial Channel Knowledge

When Eve's observations of her channel to Bob exhibit a correlation with the legitimate cascade channel, she can leverage this information in her attack. Since the observation \mathbf{z} encloses all the information Eve can exploit to deceive Bob, the attack strategy \mathbf{v} is conditionally independent of \mathbf{x} , given \mathbf{z} . We assume that \mathbf{z} is jointly Gaussian distributed with \mathbf{x} and \mathbf{y} , with joint covariance matrix

$$\mathbf{K}_{\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{bmatrix}} = \begin{bmatrix} \mathbf{K}_x & \mathbf{K}_{xy} & \mathbf{K}_{xz} \\ \mathbf{K}_{yx} & \mathbf{K}_y & \mathbf{K}_{yz} \\ \mathbf{K}_{zx} & \mathbf{K}_{zy} & \mathbf{K}_z \end{bmatrix}. \quad (27)$$

To simplify the upcoming calculation, without loss of generality, we assume that \mathbf{x} , \mathbf{y} , and \mathbf{z} are zero-mean³. In this scenario, we consider the bound in (13) since it is mathematically more tractable than (14).

From [1, Th. 2] and [13], the attack that minimizes $\mathbb{D}(p_{xv} \| p_{xy})$ when \mathbf{x} , \mathbf{y} , and \mathbf{z} are zero mean jointly Gaussian random vectors is

$$p_{v|z}^* \sim \mathcal{N}(\mathbf{G}\mathbf{z}, \mathbf{C}\mathbf{C}^H + \mathbf{K}_\omega). \quad (28)$$

Thus, the optimal attack for KL divergence minimization is

$$\mathbf{v}^* = \mathbf{G}\mathbf{z} + \mathbf{C}\boldsymbol{\omega}_c + \boldsymbol{\omega}, \quad (29)$$

where $\boldsymbol{\omega}_c \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Defining $\mathbf{P} = \mathbf{K}_{xz}^H \mathbf{K}_x^{-1} \mathbf{K}_{xz}$, we have

$$\mathbf{G} = \mathbf{K}_{xz} \mathbf{P}^\dagger, \quad (30)$$

$$\mathbf{C}\mathbf{C}^H = \mathbf{K}_{xz} \mathbf{P}^\dagger (\mathbf{I} - \mathbf{K}_z \mathbf{P}^\dagger) \mathbf{K}_{xz}^H. \quad (31)$$

We remark that $\mathbf{C}\mathbf{C}^H$ is obtained with the closed form expression in (31) if the result is a positive-semidefinite matrix, otherwise $\mathbf{C}\mathbf{C}^H$ is computed through an iterative process [1,

³Note that every agent is aware of the statistics of the random vectors, allowing the mean to be subtracted and re-added at a later stage.

Sec. IV]. We now compute $\mathbb{D}(p_{xv}||p_{xy})$ under this divergence minimizing attack.

First, we observe that z can be expressed as a linear transformation of x plus an AWGN vector $\tilde{\omega}$ with an appropriately tuned covariance matrix,

$$z = Fx + \tilde{\omega}, \quad (32)$$

where $F = K_{xz}^H K_x^{-1}$ and $K_{\tilde{\omega}} = K_z - P$. Defining $B \triangleq GF$ and $\eta = G\tilde{\omega} + C\omega_c + \omega$, (29) is rewritten as

$$v^* = Bx + \eta, \quad (33)$$

with $\eta \sim \mathcal{N}(0, K_\eta)$. Now, considering G and CC^H defined as in (30) and (31) we have

$$\begin{aligned} K_\eta &= GK_{\tilde{\omega}}G^H + CC^H + K_\omega = \\ &= K_{xz}P^\dagger K_z P^\dagger K_{xz}^H - K_{xz}P^\dagger PP^\dagger K_{xz}^H \\ &\quad + K_{xz}P^\dagger K_{xz}^H - K_{xz}P^\dagger K_z P^\dagger K_{xz}^H + K_\omega = K_\omega. \end{aligned} \quad (34)$$

When (31) is positive semidefinite, Eve gets $K_\eta = K_\omega$ and

$$\begin{aligned} D(p_{xv^*}||p_{xy}) &= \mathbb{E} \left[\log \frac{p_{v^*|x}(v|x)}{p_{y|x}(v|x)} \right] = \mathbb{E} \left[\log \frac{p_\eta(v - Bx)}{p_\omega(v - x)} \right] \\ &= \mathbb{E} \left[\log \frac{p_\eta(\eta)}{p_\omega(\eta)} \right] + \mathbb{E} \left[\log \frac{p_\omega(\eta)}{p_\omega(\eta - (I - B)x)} \right] \\ &= \mathbb{D}(p_\eta||p_\omega) + \mathbb{E} \left[((I - B)x)^H K_\omega^{-1} (I - B)x \right] \\ &= \frac{1}{2} \left[\text{tr} \left((B - I)K_x(B - I)^H K_\omega^{-1} \right) \right]. \end{aligned} \quad (35)$$

Under the same attack, we can also compute the divergence $\mathbb{D}(p_{xy}||p_{xv^*})$, in the bound (14). Following similar computation steps as (35), we get

$$\begin{aligned} \mathbb{D}(p_{xy}||p_{xv^*}) &= \mathbb{D}(p_\omega||p_\eta) + \mathbb{E} \left[\log \frac{p_\eta(\omega)}{p_\eta(\omega - (B - I)x)} \right] \\ &= \frac{1}{2} \left[\text{tr} \left((B - I)K_x(B - I)^H K_\eta^{-1} \right) \right]. \end{aligned} \quad (36)$$

Thus, in the partial channel knowledge scenario, when Eve uses the attack strategy (29), so that (34) holds, we obtain $\mathbb{D}(p_{xv^*}||p_{xy}) = \mathbb{D}(p_{xy}||p_{xv^*})$. It is important to note that v^* defined in (29) is the optimal attack strategy that minimizes $\mathbb{D}(p_{xv}||p_{xy})$, however, it is not the optimal attack strategy to minimize $\mathbb{D}(p_{xy}||p_{xv^*})$.

V. NUMERICAL RESULTS

In this Section, we validate the analytical derivations and provide numerical evidence of the effectiveness of the proposed KL divergence minimizing attack. As detection test, Bob uses the generalized likelihood ratio test (GLRT), which is done by classifying the observed channel r (see Fig. 2) legitimate, for a fixed threshold τ , if

$$\Lambda(r) = \log \frac{\max_v p_{v|x}(v|x)}{p_{y|x}(r|x)} \leq \tau. \quad (37)$$

To simulate the scenario, we consider correlated Rayleigh fading channels. Specifically, we assume that the entries of H and G are random, independent within each matrix and

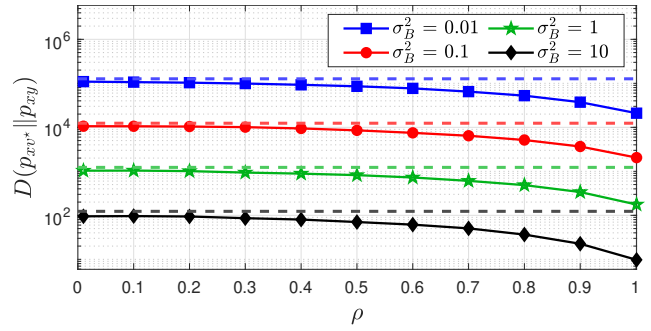


Fig. 3. KL divergence vs ρ , one curve for every value of σ_B^2 , with $N_A = N_B = N_E = 5$ antennas, and $\gamma = \pi/3$. Dashed lines represent (23) for each value of σ_B^2 .

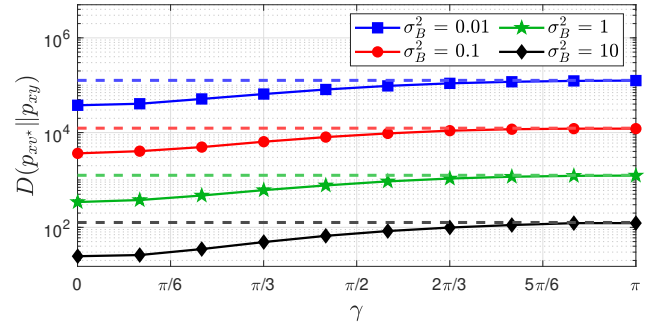


Fig. 4. KL divergence vs γ , one curve for every value of σ_B^2 , with $N_A = N_B = N_E = 5$ antennas, and $\rho = 0.8$. Dashed lines represent (23) for each value of σ_B^2 .

between the two matrices, and Gaussian distributed with zero mean and unitary variance. Moreover, we assume that $N_A = N_E$ and the observation Z of Eve is correlated with the legitimate cascade channel as follows

$$Z = \rho H \Phi^* G + \sqrt{1 - \rho^2} D, \quad (38)$$

with $\rho \in [0, 1]$ representing the correlation factor. Furthermore, D entries are independent of each other, independent of those of both G and H , and follow a Gaussian distribution with zero mean and unitary variance. We set $N = 100$ IRS elements so that the approximation of the distribution of x , y , and z as Gaussian is meaningful.

Fig. 3 shows $\mathbb{D}(p_{xv^*}||p_{xy})$ vs ρ , when Eve uses the optimal attack strategy (29). The dashed lines represent the KL divergence (23), obtained with attack (22), which is optimal for $\rho = 0$. We show results for different values of σ_B^2 , with $N_A = N_B = N_E = 5$ antennas. We note that as $\rho \rightarrow 0$, the solid curves approached the dashed lines. Thus, the optimal attack (29) for the partial channel knowledge scenario converges, as ρ goes to zero, to the optimal attack (22) in the scenario without channel knowledge. Moreover, the KL divergence increases significantly with σ_B^2 .

Fig. 4 shows $\mathbb{D}(p_{xv^*}||p_{xy})$ vs γ for $\rho = 0.8$, and $N_A = N_B = N_E = 5$. We see that the KL divergence grows with γ , as expected, since the higher the γ , the more the IRS configurations chosen by Bob will deviate from the

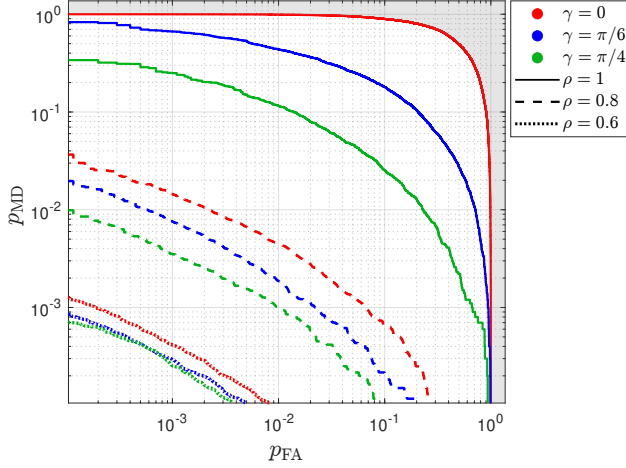


Fig. 5. DET curves for the GLRT for different value of ρ , with $N_A = N_E = 5$, $N_B = 2$ antennas, and $\sigma_B^2 = N/10$. The area corresponding to $p_{MD} \geq 1 - p_{FA}$ is shaded, and its edge represents the trivial limit case in which the decision is made tossing a biased coin.

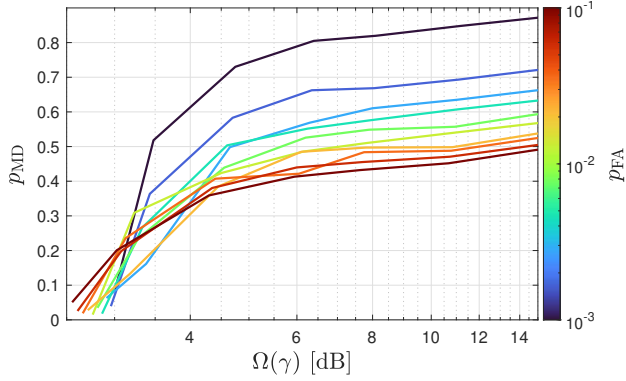


Fig. 6. p_{MD} as a function of $\Omega(\gamma)$ for various values of p_{FA} with $N_A = N_E = N_B = 1$ antenna, $\rho = 0.6$, and $\sigma_B^2 = N/2$.

communication-optimal one (3). Moreover, the KL divergence saturates to (23) when $\gamma \rightarrow \pi$.

Fig. 5 shows the detection error tradeoff (DET) curves for different value of ρ , with $N_A = N_B = N_E = 5$ antennas, and $\sigma_B^2 = N/10$. As ρ increases, the DET curves move towards the edge of the shaded area, which represents the trivial limit case when the decision is taken randomly, without using the signal. Thus, the FA probability for a given MD probability increases when the correlation between Eve's observation and the actual cascade Alice-IRS-Bob channel is higher.

The average SNR in the considered scenario is an indicator of communication performance and is defined as

$$\Omega(\gamma) = (1 - p_{FA}) \mathbb{E} \left[\frac{1}{\sigma_B^2} \left| \sum_{n=1}^N H_n G_n e^{j\theta_n} \right|^2 \right] \approx (1 - p_{FA}) \frac{N}{\sigma_B^2} \left((N-1) \frac{\pi^2 \sin(\gamma)^2}{16\gamma^2} + 1 \right), \quad (39)$$

where the approximation holds asymptotically ($N \rightarrow \infty$ by the central limit theorem). Fig. 6 shows the p_{MD} as a function

of the average SNR when Eve uses the optimal attack strategy (29), for $p_{FA} \in [10^{-3}, 10^{-1}]$ and considering a single antenna scenario, i.e., $N_A = N_E = N_B = 1$. From Fig. 6 we see the trade-off between p_{MD} and Ω . Specifically, Ω grows with γ . However, achieving a significantly higher Ω comes at the cost of a higher p_{MD} , while choosing a lower γ setting corresponds to a high p_{MD} and thus greater vulnerability to attack.

VI. CONCLUSIONS

In this paper, we have investigated new attacks within the context of a CR-PLA, building upon the findings of [1], i.e., on performance bounds (for FA and MD probabilities) for conventional PLA. We have provided an analysis of CR-PLA with IRS, which falls under the assumptions of [1] when a large number of IRS elements is considered. We have derived the optimal attack strategy when the attacker has either partial or no knowledge of the cascade Alice-IRS-Bob channel. We have assessed the performance of the proposed attack under various correlations between Eve's observation and the legitimate cascade channel, and have presented results on the p_{MD} vs p_{FA} tradeoff for GLRT, showing the effectiveness of the proposed attack.

REFERENCES

- [1] A. Ferrante, N. Laurenti, C. Masiero, M. Pavon, and S. Tomasin, "On the error region for channel estimation-based physical layer authentication over Rayleigh fading," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 5, pp. 941–952, 2015.
- [2] G. J. Simmons, "Authentication theory/coding theory," in *Advances in Cryptology*, G. R. Blakley and D. Chaum, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 1985, pp. 411–431.
- [3] P. Baracca, N. Laurenti, and S. Tomasin, "Physical layer authentication over MIMO fading wiretap channels," *IEEE Trans. Wireless Commun.*, vol. 11, no. 7, pp. 2564–2573, 2012.
- [4] W. Hou, X. Wang, J.-Y. Chouinard, and A. Refaey, "Physical layer authentication for mobile systems with time-varying carrier frequency offsets," *IEEE Transactions on Communications*, vol. 62, no. 5, pp. 1658–1667, 2014.
- [5] R. Diamant, P. Casari, and S. Tomasin, "Cooperative authentication in underwater acoustic sensor networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 2, pp. 954–968, 2019.
- [6] U. M. Maurer, "Authentication theory and hypothesis testing," *IEEE Transactions on Information Theory*, vol. 46, no. 4, p. 1350–1356, 2000.
- [7] H. Fang, X. Wang, and L. Hanzo, "Learning-aided physical layer authentication as an intelligent process," *IEEE Transactions on Communications*, vol. 67, no. 3, pp. 2260–2273, 2019.
- [8] E. Jorswieck, S. Tomasin, and A. Sezgin, "Broadcasting into the uncertainty: Authentication and confidentiality by physical-layer processing," *Proceedings of the IEEE*, vol. 103, no. 10, pp. 1702–1724, 2015.
- [9] N. Xie, Z. Li, and H. Tan, "A survey of physical-layer authentication in wireless communications," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 1, pp. 282–310, 2021.
- [10] S. Tomasin, H. Zhang, A. Chorti, and H. V. Poor, "Challenge-response physical layer authentication over partially controllable channels," *IEEE Communications Magazine*, vol. 60, no. 12, pp. 138–144, 2022.
- [11] F. Mazzo, S. Tomasin, H. Zhang, A. Chorti, and H. V. Poor, "Physical-layer challenge-response authentication for drone networks," in *Proc. IEEE Global Commun. Conference (GLOBECOM)*, 2023.
- [12] Ö. Özdogan, E. Björnson, and E. G. Larsson, "Using intelligent reflecting surfaces for rank improvement in MIMO communications," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 9160–9164.
- [13] L. Crosara, F. Ardizzone, S. Tomasin, and N. Laurenti, "Worst-case spoofing attack and robust countermeasure in satellite navigation systems," *IEEE Transactions on Information Forensics and Security*, vol. 19, pp. 2039–2050, 2024.