

Adversarial Bandit Approach for RIS-Aided OFDM Communication

Messaoud Ahmed Ouameur

Universite du Quebec a Trois-Rivieres

Dương Tuấn Anh Lê

VNUHCM-US: University of Science

Daniel Massicotte (✉ daniel.massicotte@uqtr.ca)

Université du Québec à Trois-Rivières <https://orcid.org/0000-0002-7807-7919>

Gwanggil Jeon

Incheon National University

Felipe Augusto Pereira de Figueiredo

National Telecommunications Institute

Research

Keywords: reconfigurable intelligent surfaces, reflection beamforming prediction, deep learning, machine learning, sixth generation (6G) wireless systems, adversarial bandit, exponential-weight algorithm for exploration and exploitation, follow the perturbed leader (FPL)

Posted Date: February 2nd, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1244428/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Adversarial Bandit Approach for RIS-Aided OFDM Communication

Messaoud Ahmed Ouameur¹, Lê Dương Tuấn Anh², Daniel Massicotte^{1*}, Gwanggil Jeon³ and Felipe Augusto Pereira de Figueiredo⁴

* Correspondence: daniel.massicotte@uqtr.ca

¹ Department of Electrical and Computer Engineering, Université du Québec à Trois-Rivières, Québec, Canada
Full list of author information is available at the end of the article

Abstract— To assist sixth-generation wireless systems in the management of a wide variety of services, ranging from mission-critical services to safety critical tasks, key physical layer technologies such as reconfigurable intelligent surfaces (RISs) are proposed. Even though RISs are already used in various scenarios to enable the implementation of smart radio environments, they still face challenges with regards to real time operation. Specifically, RISs typically need costly, high dimensional channel estimation with offline exhaustive search, requiring prohibitive hardware complexity or online exhaustive beam-training that incurs high training overhead. While in its infant stage, the application of deep learning (DL) tools shows promise in enabling feasible solutions. In this paper, we propose two low-training overhead and energy-efficient adversarial bandit-based schemes with outstanding performance gains when compared to DL-based reflection beamforming reference methods. The resulting deep learning models are discussed using state-of-the-art model quality prediction trends.

Index Terms— reconfigurable intelligent surfaces, reflection beamforming prediction, deep learning, machine learning, sixth generation (6G) wireless systems, adversarial bandit, exponential-weight algorithm for exploration and exploitation, follow the perturbed leader (FPL).

I. Introduction

Sixth-generation (6G) wireless systems are expected to enable greater levels of autonomy, improve human-machine interfacing, and achieve deep connectivity in more diverse environments. To assist 6G in managing a wide variety of services, ranging from mission critical services (e.g. autonomous driving) to safety critical tasks (e.g. remote surgery), key enabling physical layer technologies (PHY) such as ultra-massive multiple-input multiple-output (MIMO) systems, millimeter wave and Tera-Hertz communications, and reconfigurable intelligent surfaces (RISs), need to be carefully designed [1]. Unfortunately, current network design practices conform to the hypothesis that regards the wireless environment between communicating devices to be unmodified and can be only overcome through the design of advanced transmission and reception schemes. Breaking free from such a hypothesis by programming the environment is expected to enable major performance gains. As such, RIS-aided communication has received increasing interest from the research community due to its potential in extending the coverage, enhancing link quality and capacity, and mitigating interference and security breaches [2]. RISs enable the reconfiguration of the wireless propagation environment by intelligently controlling the signal reflections via its massive low-cost elements. By jointly adapting the reflected signal amplitude and/or phase shift at each RIS element based on the wireless channels, the signals

reflected by the RIS can be constructively combined at the intended receiver. Unlike traditional active relaying/beamforming techniques, RIS is designed to be totally or nearly passive, thus enjoying lower hardware cost and energy consumption [1]. So far, RIS has been adopted in various scenarios. In [3], the error performance of an RIS-aided single-input single-output (SISO) system is examined; meanwhile, RISs are also used for multi-user systems to maximize the signal-interference-noise ratio [5] or to enhance energy efficiency [4]. Unfortunately, due to the additional channel links between the RIS and its associated transmitter and intended receivers, the large gain is achieved at the expense of increased overhead for channel estimation. Early works focus on the design of reflection beamforming coefficients under the assumption of perfect channel state information [6], which helps in deriving the system performance bounds, but the underlying optimal techniques are unfortunately algorithm-deficient. Obtaining this channel knowledge, in practice, may require large and possibly prohibitive training overhead, which represents the main challenge for real-time RIS operation. In addition, implementing the RIS using discrete (and possibly non-accurate) phase shifters makes it difficult to analytically model such behavior in a tractable manner, making the overall end-to-end model-deficient. Under such deficiencies, machine learning is introduced and has started to be extensively used to enhance the implementation of various components within the 5G radio access network (RAN) [7]. In addition to the smart radio environment concept, embracing a vision wherein 6G is designed in a way that ML could modify parts of the physical (PHY) and medium access control (MAC) layers is proposed [7]. Deep Learning (DL) has also been used for devising computationally efficient approaches for physical layer communication receiver modules. Under the supervised learning approach, the authors in [8] present a DL framework for MIMO symbol detection. It has been able to achieve near-optimal detection performance with an even faster real-time implementation. A recurrent neural network (RNN)-based detection scheme is introduced in [9] for MIMO orthogonal frequency division multiplexing (OFDM) systems and is shown to outperform traditional detection techniques under channel impairments and hardware non-linearities. Convolutional neural network (CNN)-based supervised learning techniques can also be utilized for channel estimation problems, providing improved generalization abilities and robustness to channel alterations [10]. A DL-based beam prediction method was proposed for distributed mmWave MIMO systems to cope with highly mobile users with negligible training overhead and high data rate gains [11]. The authors in [12] present a novel RIS hardware architecture along with two solutions based on compressive sensing and deep reinforcement learning with very negligible training overhead. Deep reinforcement learning (DRL) has also been applied for designing efficient spectrum access [13] and scheduling strategies [14] for cellular networks. Automatic cell-sectorization for cellular network coverage maximization is another area where DRL has shown tremendous potential [15].

In this paper we propose two efficient reinforcement learning based schemes where the main contributions are as follows:

- We propose an adversarial bandit approach based on exponential-weight algorithm for exploration and exploitation (EXP3). To show the merits of the proposed scheme, we conduct extensive

simulation using the publicly available accurate ray tracing based DeepMIMO dataset [16], with the 'O1' scenario.

- To improve upon the computational complexity, the Follow the Perturbed Leader (FPL) scheme is discussed.
- To compare the quality of the state-action deep neural network models used with the reference methods in [12] and with the proposed ones (EXP3 and FPL), we leverage state of the art techniques such as the power law (PL) exponents [17].

The paper is organized as follows: The system model and problem formulation are presented under method section II. Section II discusses also the proposed adversarial bandit approaches. Section III is devoted to discuss the results in terms of achievable rate and energy efficiency while considering a low complexity alternative using a FPL algorithm. The associated DL models' quality is also analyzed using PL exponents. Finally, the conclusions are made and future research directions are outlined in Section V.

II. Methods

The independent and identically distributed Rayleigh fading channel is not physically present when using RIS with a rectangular arrangement. Therefore, an alternative physically feasible model for evaluating RIS-aided communications is required [18]. To enable practical implementations of RIS-aided communication systems, new path loss models [18], [2], and open-source channel models [2], [16] have been developed. As such to reproduce the results and perform a fair comparison, we will adopt the system and channel model in [12].¹

A. System model

As depicted in Figure 1, transmitter-receiver communication is aided by an RIS having M reconfigurable elements. For the sake of simplicity, we assume that both the transmitter and receiver are equipped with a single antenna. For generalization, one can adopt the signal model from [2]. An OFDM-based transmission with K subcarriers is adopted. The direct channel per subcarrier k between the transmitter and the receiver is denoted by $h_{\text{TR},k} \in \mathbb{C}$ whereas links via the RIS are represented by $M \times 1$ complex valued vectors $\mathbf{h}_{\text{T},k}, \mathbf{h}_{\text{R},k} \in \mathbb{C}^{M \times 1}$. By neglecting the direct path², the received signal can be written as

$$y_k = \mathbf{h}_{\text{R},k}^T \mathbf{\Psi}_k \mathbf{h}_{\text{T},k} s_k + n_k \quad (1)$$

Where $\mathbf{\Psi}_k \in \mathbb{C}^{M \times M}$ is the RIS interaction diagonal matrix, s_k and n_k are the transmitted symbol per subcarrier k and the receive noise with zero mean and variance of σ_n . With P_{T} being the total transmit

¹ When using DL tools, it is hard to evaluate the merits and the performance of the proposed methods in comparison with reference methods unless similar models and datasets are used. Otherwise, any performance gain may be attributed to system and channel model and dataset differences.

² The benefit of RIS is mostly harnessed when the direct path is blocked or simply very weak. Such an assumption helps in simplifying the analysis of the algorithm [1] [2] [12] [16].

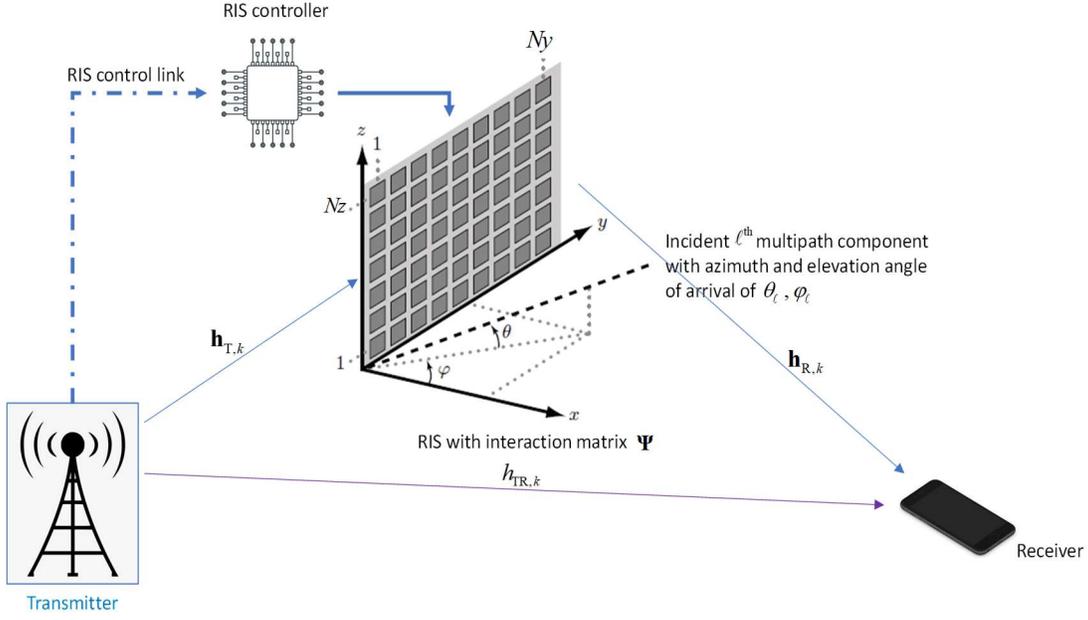


Figure 1. The system model in which transmitter-receiver communication is aided by a RIS having an $M \times M$ interaction matrix where $M = N_y \cdot N_z$

power, the follow power constraint per subcarrier is enforced $E(|s_k|^2) = P_T/K$. Herein $(\cdot)^T$ and $E(\cdot)$ denote the transpose and the expectation operations respectively. If we re-arrange the diagonal elements of the interaction matrix Ψ_k in an $M \times 1$ column vector ψ_k , we refer to it as the reflection beamforming (BF) vector, such that $\Psi_k = \text{diag}(\psi_k)$, equation (1), can also be expressed in more convenient way as

$$y_k = (\mathbf{h}_{R,k} \odot \mathbf{h}_{T,k})^T \psi_k s_k + n_k \quad (2)$$

where \odot denotes the Hadamard product. Imposing few practical implementation constraints of a nearly passive RIS where the phase shifters apply the same phase shift over all subcarriers, the m -th element of ψ_k is modeled as $[\psi_k]_m = e^{j\phi_m}$.

B. Channel model and RIS design objective

As for the channel model, a wideband geometric channel model for $\mathbf{h}_{T,k}$ and $\mathbf{h}_{R,k}$ is used [2] [12] [18]. Therefore, $\mathbf{h}_{T,k}$ and $\mathbf{h}_{R,k}$ is expressed as a function of the azimuth and elevation angles of arrival/departure of the ℓ^{th} path from a total of L paths such that the array vector of the RIS is defined as $\mathbf{a}(\theta_\ell, \varphi_\ell) \in \mathbb{C}^{M \times 1}$ where $\theta_\ell \in [0 \ 2\pi)$ and $\varphi_\ell \in [0 \ \pi)$ (see Figure 1). For the sake of brevity, we refer the reader to [2] and [12] for detailed modeling.

The RIS design objective is therefore to find out the reflection BF vector ψ_k that maximizes the achievable rate at the receiver

$$R = \frac{1}{K} \sum_{k=1}^K \log_2 \left(1 + \rho \left| \left(\mathbf{h}_{R,k} \odot \mathbf{h}_{T,k} \right)^T \boldsymbol{\Psi}_k \right|^2 \right) \quad (3)$$

where the signal to noise ratio is $\rho = P_T / K \sigma_n$. The maximization is done over a discrete pre-defined codebook P due to the fact that a practical radio frequency (RF) phase shifter uses quantized phase values. Unfortunately, maximizing (3) entails an exhaustive search over the codebook P . Fortunately, the authors in [12] have proposed a novel hardware architecture along with a compressive sensing and DL-based framework to tackle the issue with low training overhead. However, there is still a large room for improvement as we will discuss throughout this paper.

C. Proposed algorithm using Adversarial bandit approach via exponential-weight algorithm for exploration and exploitation

The authors in [12] use a DL-based approach to predict the reflection BF vector. Over a channel coherence block size S , the RIS sends two pilots to estimate a sampled channel vector $\bar{\mathbf{h}}(s) = \text{vec}([\bar{\mathbf{h}}_1(s), \bar{\mathbf{h}}_2(s), \dots, \bar{\mathbf{h}}_K(s)])$ where $\bar{\mathbf{h}}_k(s) \in \mathbb{C}^{\bar{M} \times 1}$ denotes the sampled combined channel vector, $\mathbf{h}_k = \mathbf{h}_{R,k} \odot \mathbf{h}_{T,k}$, for the k -th subcarrier at s -th channel coherent block using a fraction number of the RIS elements $\bar{M} \ll M$ that are assumed to be active elements (i.e. equipped with full RF and baseband processing chain for an effective uplink and downlink channel estimation). During beam training, the RIS is configured using one reflection beam $\boldsymbol{\Psi}$ (notice that the subscript k is removed because one reflection BF vector is available for all subcarriers) from the codebook P . Then, a dataset is contracted out of the tuples $\Upsilon \leftarrow (\bar{\mathbf{h}}(s), \mathbf{r}(s))$ where $\mathbf{r}(s) = [R_1(s), R_2(s), \dots, R_N(s)]^T$ and $R_n(s)$ is the measured rate using the n -th codebook (N is the cardinality of the codebook P). Finally, a deep neural network is trained using the dataset Υ .

D. Adversarial bandit approach via exponential-weight algorithm for exploration and exploitation

Despite the novel architecture that suggests the use of a few active elements to sample the uplink and downlink channel vectors, the proposed algorithm can be substantially improved. As such, we propose an approach based on adversarial bandit wherein instead of spanning equally every element of the codebook P , we adopt a scheme that favors the more likely optimal beams. Therefore, the dataset Υ will have more useful data to train with. Table 1 shows the proposed adversarial bandit based on exponential-weight algorithm for exploration and exploitation (EXP3) [21]. The adversarial bandit is a variant of the multi-armed bandit problem where a fixed limited set of resources (phase shifters) must be assigned among alternative choices (reflection beamforming) in a way that maximizes their expected gain (achievable rate), when the properties of each choice are only partly known at the time of assignment and may become better comprehended as time passes. This is one of the strongest generalizations of the bandit problem as it disregards all assumptions of the distribution. In its basic form [27], EXP3 chooses a reflection beamforming vector $\boldsymbol{\Psi}_n$ (steps 4 and 5 in Table 1) from the codebook

Table 1. Adversarial bandit based scheme for reflection beamforming vector perdition

Input: reflection beamforming codebook	
Initialize: $\gamma \in (0,1]$ and weights $w_n(t=1) = 1$ for $n = 1, \dots, N$	
Phase I: Learning phase	
1	For $s = 1$ to S do (span over S channel coherent blocks)
2	RIS receives two pilots to estimate $\bar{\mathbf{h}}(s)$
3	For $t = 1$ to T do (go over T iterations)
4	Compute the probability of the n -th codebook $p_n(t) = (1-\gamma) \frac{w_n(t)}{\sum_{c=1}^N w_c(t)} + \frac{\gamma}{N}$
5	Draw n_t randomly according to the probabilities $p_1(t), p_2(t), \dots, p_N(t)$ and RIS uses $\boldsymbol{\Psi}_{n_t}$
6	RIS receives $R_{n_t}(s)$
7	Compute the reward $r_{n_t} \in [0,1]$
8	For $c = 1$ to N do (update the weights)
9	$\hat{r}_c(t+1) = \begin{cases} r_c(t)/p_c(t) & \text{if } c = n_t \\ 0 & \text{otherwise} \end{cases}$
10	$w_n(t+1) = w_n(t) e^{\gamma \hat{r}_c(t)/N}$
11	Store the new entry in the learning data set $\Upsilon \leftarrow (\bar{\mathbf{h}}(s), \mathbf{p}(s))$ where $\mathbf{p}(s) = [p_1(T), p_2(T), \dots, p_N(T)]^T$
12	Train the DL model using the learning dataset Υ
Phase II: Prediction phase	
13	While True do (for every channel coherent block s')
14	RIS receives two pilots to estimate $\bar{\mathbf{h}}(s')$
15	Predict the probability vector $\mathbf{p}(s')$
16	RIS uses $\boldsymbol{\Psi}_{n^*}$ where $n^* = \arg \max_n [\mathbf{p}(s')]_n$

P at random with probability $(1-\gamma)$ where it prefers choices with higher weights (exploit), or it selects with probability γ to uniformly randomly explore. After receiving the rewards (steps 6 and 7), the weights are updated (steps 9 and 10). The exponential growth significantly increases the weight of good reflection beamforming vectors.

The key advantage of the proposed scheme is that, instead of using the rates as the deep neural network outputs, we use the pull-probability vector $\mathbf{p}(s) = [p_1(T), p_2(T), \dots, p_N(T)]^T$ computed at step 4 using the updated weights which are in turn computed using the normalized reward (step 7). The

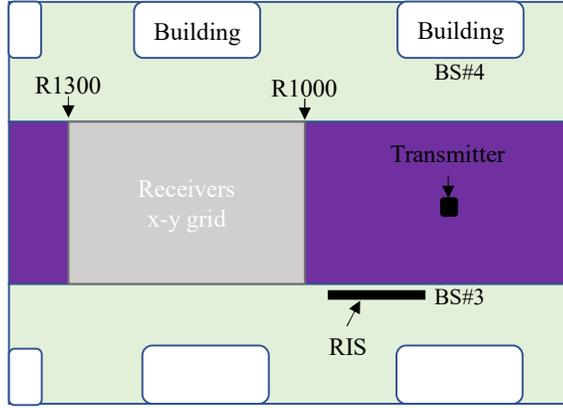


Figure 2. The ray tracing scenario ‘O1’ from [26]. The BS#3 is selected to be the RIS while the transmitter is fixed at the position of raw R850 and column 90. The receiver can be positioned at any 54300 points within the x-y grid between raw R1000 and R1300. These points constitute the dataset which is split into 80% training set and 20% test set.

normalized reward r_{n_t} is computed using the received rate $R_{n_t}(s)$ as $r_{n_t} = R_{n_t}(s)/R_{\max}$, where R_{\max} is the maximum achievable rate or a large number to make sure that $r_{n_t} \in [0,1]$.

III. RESULTS AND DISCUSSIONS

The proposed EXP3-based learning scheme is evaluated using the outdoor ray tracing scenario O1 from the deep-MIMO dataset that is publically available at [16]. For the sake of facilitating the comparison, a similar setup is used in [12] as well (see Figure 2). The results herein are also validated using channel data generated using SimRIS tool [2]. The adopted RIS employs a uniform planar array (UPA) with 16-by-16 ($M = 256$) antenna elements with 3 dBi gain at the 28 GHz mmWave setup. The transmit power is set to 10 dBW while the receiver’s noise figure is 5 dB. The codebook P is constructed using a 2D discrete Fourier transform (DFT) matrix.

The number of subcarriers involved in $\bar{\mathbf{h}}(s)$ is ($\bar{K} = 64$) \ll ($K = 512$), which sets the input of the DL model equal to $2\bar{K}\bar{M}$. The sampled channel vector is normalized prior to the training phase. The DL models consists of four layers similar to the one used in [12] where the number of the nodes in the hidden layers is $(2\bar{K}\bar{M}, 4M, 4M, M)$. The regular training and optimization parameters are: batch size set to 500 samples, drop-out rate is 0.5 and L_2 regularization factor is 0.0001. Of course, we do not attempt to optimize the DL model but we will discuss its quality using the state-of-the-art techniques such as the power law exponents [17] in section IV.

Figure 3 shows the achievable rate as a function of the number of training samples. The proposed EXP3-based scheme requires substantially less data as compared to DL reflection beamforming technique [12], owing to the optimal selection of the dataset which stresses that less likely reflection beams are given lower probability which excludes them during the exploitation phase of the EXP3 algorithm (Table 1, Step 5). The reference DL reflection beamforming requires more active elements \bar{M}

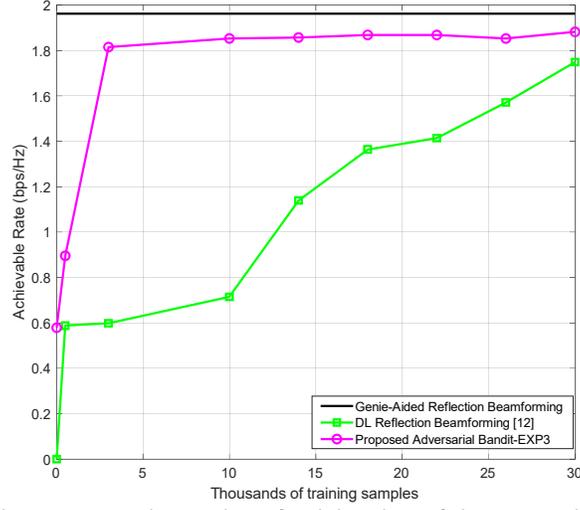


Figure 3. The achievable rate versus the number of training data of the proposed EXP3-based scheme in comparison with the reference DL reflection beamforming [12] and the reference genie-aided method (that assumes perfect knowledge of the channel) where $\bar{M} = 4$

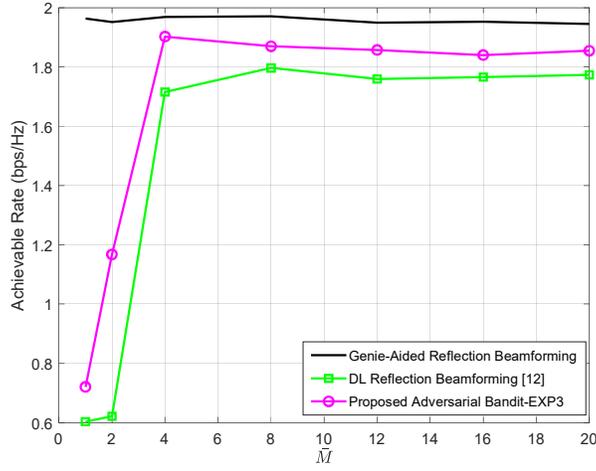


Figure 4. The achievable rate versus the number of active elements of the proposed EXP3-based scheme in comparison with the reference DL reflection beamforming [12] and the reference genie-aided method (that assumes perfect knowledge of the channel) where the number of training data is set to 30K

to sustain competitive performance as shown in Figure 4 where EXP3-based learning schemes achieves 96% of the optimal achievable rate compared to 88% using the reference method in [12]. However, this will come at the expense of higher power consumption. Nevertheless, it seems that as far as the number of active elements is higher than 4, all methods are showing close to the performance of the genie-aided method.

We reformulate the energy efficiency as $\eta = W \times R / P_c$ measured in Mbit/J, where W is the transmission bandwidth and P_c is the RIS power consumption which can be broken down to

$$P_c = MP_{PS} + \bar{M} (P_{LNA} + P_{RF} + 2FOM_W f_{FS} 2^b) \quad (4)$$

where the term $2FOM_W f_{FS} 2^b$ is the power consumption of a b -bits ADC with f_{FS} being the Nyquist sampling frequency and FOM_W is the Walden's figure of merit [19]. P_{PS} , P_{LNA} and P_{RF} are, respectively,

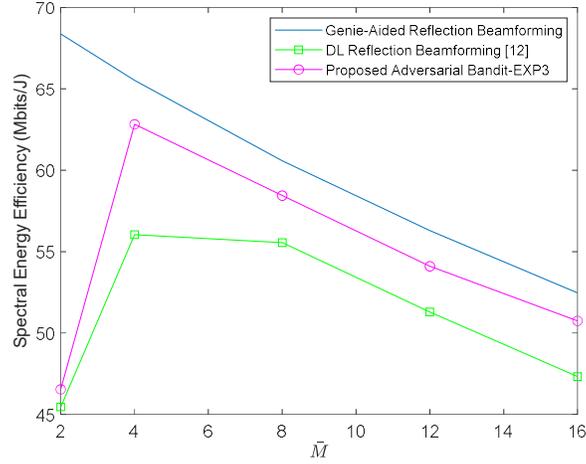


Figure 5. The energy efficiency η versus the number of active elements \bar{M} of the proposed EXP3-based scheme in comparison with the reference DL reflection beamforming [12] and the reference genie-aided method where the number of training data is set to 30K

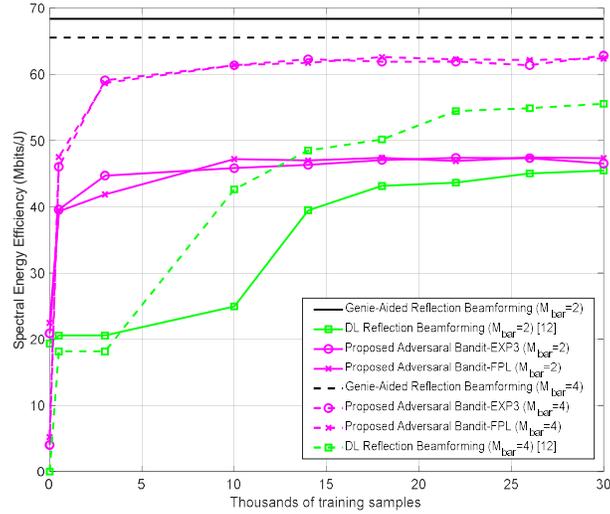


Figure 6. The spectral energy efficiency versus the number of training samples and the number of active elements \bar{M} for the proposed FPL-based scheme.

the power consumptions of the phase-shifter in the passive RF path, and the low-noise amplifier (LNA) and the rest of the RF chain along the active paths. As per the state-of-the-art RF parts' specifications, these variables are set to $P_{PS} = 10$ mW, $P_{LNA} = 20$ mW, $P_{RF} = 40$ mW and the baseband processing power of 200 mW is assumed. Assuming similar values like the ones in [12] and [20], $FOM_w = 46.1$ fJ/conversion at $W = 100$ Mhz and $b = 4$ bits. As such, Figure 5 depicts the energy efficiency η as a function of the number of active elements \bar{M} . Like the reference DL reflection BF [12], the proposed method shows optimal but higher energy efficiency performance using 4 active elements only.

In light of these results, the EXP3-based adversarial bandit method demonstrates outstanding performance gains compared to other state-of-the-art methods. So far, the adopted deep neural network

Table 2. Follow the perturbed leader (FPL) scheme for reflection beamforming perdition

Input: reflection beamforming codebook	
Initialize: γ real value	
Rewards $\hat{r}_n(t=1) = 0$ for $n = 1, \dots, N$	
Phase I: Learning phase	
1	For $s = 1$ to S do (span over S channel coherent blocks)
2	RIS receives two pilots to estimate $\bar{\mathbf{h}}(s)$
3	For $t = 1$ to T do (go over T iterations)
4	Generate a random noise from an exponential distribution for each arm $n : z_n(t) \sim \exp(\gamma)$
5	Draw an arm n_t where $n_t = \arg \max_n [\hat{r}_n(t) + z_n(t)]$ RIS uses $\boldsymbol{\Psi}_{n_t}$
6	RIS receives $R_{n_t}(s)$
7	Compute the instantaneous reward $r_{n_t} \in [0, 1]$ based on $R_{n_t}(s)$
8	For $c = 1$ to N do (update the rewards)
9	$\hat{r}_c(t+1) = \begin{cases} \hat{r}_n(t) + r_c(t) & \text{if } c = n_t \\ 0 & \text{otherwise} \end{cases}$
10	Store the new entry in the learning data set $\Upsilon \leftarrow (\bar{\mathbf{h}}(s), \hat{\mathbf{r}}(s))$ where $\hat{\mathbf{r}}(s) = [\hat{r}_1(T), \hat{r}_2(T), \dots, \hat{r}_N(T)]^T$
11	Train the DL model using the learning dataset Υ
Phase II: Prediction phase	
12	While True do (for every channel coherent block s')
13	RIS receives two pilots to estimate $\bar{\mathbf{h}}(s')$
14	Predict the probability vector $\hat{\mathbf{r}}(s')$
15	RIS uses $\boldsymbol{\Psi}_{n^*}$ where $n^* = \arg \max_n [\hat{\mathbf{r}}(s')]_n$

architecture is similar to the one used in [12]. The reason being that one would be keen to see the effect of using a new learning scheme rather than proposing a new DL model. The other reason, which we will discuss in the next section, is that one will also be interested to compare the quality of the two networks trained using $\Upsilon \leftarrow (\bar{\mathbf{h}}(s), \mathbf{r}(s))$ for [12] and $\Upsilon \leftarrow (\bar{\mathbf{h}}(s), \mathbf{p}(s))$ in the proposed method. However, let's first introduce another computationally efficient adversarial bandit-based scheme that uses the Follow the Perturbed Leader (FPL) algorithm.

A. Improving and evaluation of the quality of proposed approaches

Even if the EXP3 algorithm has efficient theoretical guarantees, it is computationally expensive due to the calculation of the exponential terms [21]. The FPL algorithm is then introduced to alleviate the burden by following the reflection beam that has the best performance while adding exponential noise to it to provide exploration [22]. Even though the baseline FPL algorithm does not have appreciated

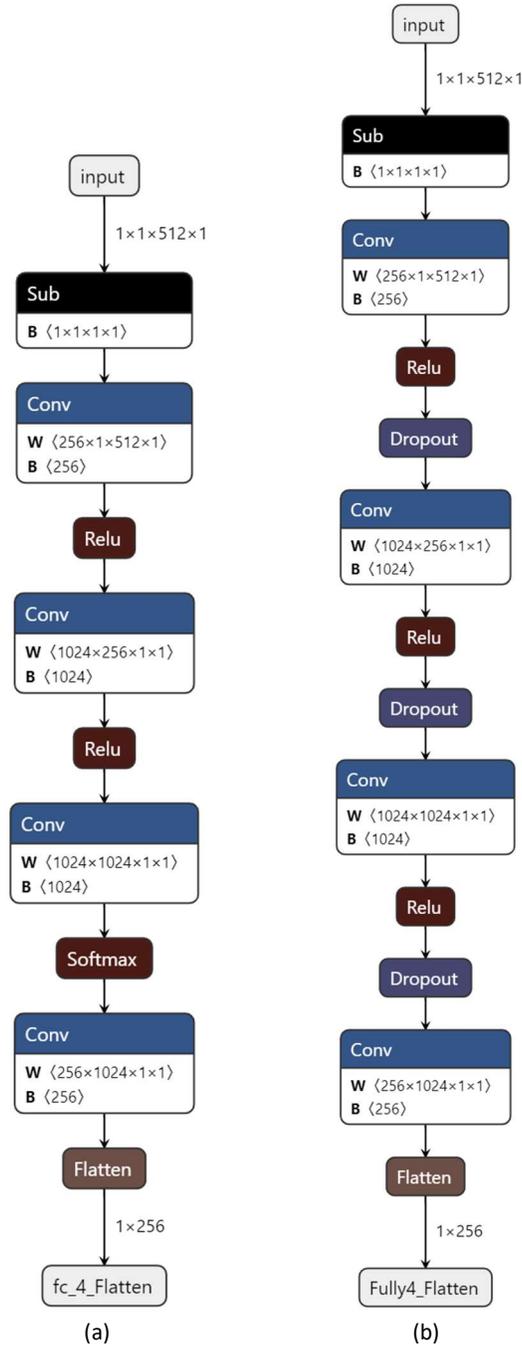


Figure 7. The DL models used with (a) EXP3/FTPL algorithm and with (b) the reference method [12] generated using Netron.

theoretical guarantees, it is worth evaluating its performance in the scope of the current RIS reflection beamforming prediction. Table 2 shows the FPL algorithm where the exponential noise, which can be computed offline, is added in step 4 to provide exploration.

Figure 6 shows that the FPL algorithm provides similar performances to the EXP3 algorithm at the expense of less “explainability” information, such as the pull-probabilities and weights inherent in EXP3. However, how one can decide which algorithm is better beyond just comparing the achievable

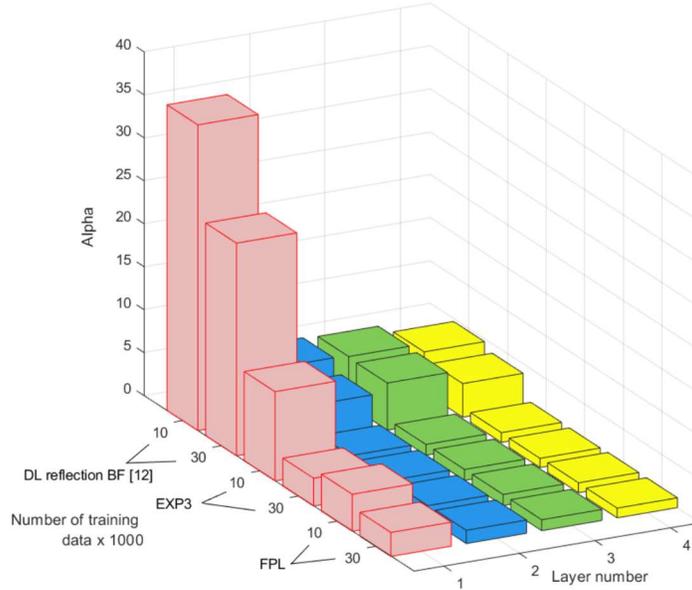


Figure 8. The power law (PL) exponent α to compare the quality of the models used in the reference method [12], EXP3 and FLP algorithms.

rates (accuracy)? Even if all algorithms have different approaches to build the training dataset \mathcal{Y} , they all share a similar model. Figure 7³ depicts the DL models used with the EXP3/FPL algorithms and the reference method [12]. The slight differences are in using the dropout layers to improve the regularization of the reference method and the use of the softmax activation for the model used with EXP3 to generate the pull-probabilities. Nevertheless, in the end, these models are considered as black-boxes that need to be compared.

It is beyond the scope of this paper to dig into explainability of DL models which can be found in [23]. We will rather use state-of-the-art tools from [17] and mainly power law (PL) exponents to compare the quality of the DL models. Figure 8, depicts the PL exponents for the 4 layers. Indeed, we expect that a poorly trained model will lack good (i.e., small exponents α) PL behavior in some layers, whereas the EXP3 has, on average, smaller α values than the reference method, with all $\alpha \leq 6$ and with smaller mean/median α . It also has far fewer unusually large outlying α values than the reference method. The model used with FPL algorithm is rather showing the best training quality at the expense of less theoretical guarantees. The exponent values are obtained using the WeightWatcher tool from [24]. For future investigation, this should also be contrasted with the behavior displayed by scale-dependent metrics such as the Frobenius norm and the Spectral norm [17].

IV. Conclusion

RIS-aided communication has received increasing interest from the research community, with discussions not just about its unprecedented potential but also about the stumbling blocks with regards to feasible real-time operation. Among others, channel estimation overhead is regarded as serious issue,

³ Netron is used to generate Figure 7. It is a visualizer for neural network, deep learning, and machine learning models. It can be acquired from <https://www.electronjs.org/apps/netron>.

which makes the adoption of DL tools an attractive alternative to solve the problem. As such, we have discussed two adversarial bandit-based schemes that provide substantial spectral and energy efficiency gains. We have also discussed the associated DL models' quality using the PL exponents to show the training quality using the dataset generated from the proposed schemes. Our work contributes to shed light on the potential improvements that can be made in exploring the interplay between ML and RISs. For future research, one could investigate the proposed schemes under different channel and system parameters while adopting meta-learning approach [25]-[26] to improve the online training performance. Different DL models can also be used along with these schemes wherein explainability shall be given a considerable attention [23] to improve the trustworthiness of the DL-enabled solutions. Last but not least, at the hardware level, one can use low power root-mean-square and envelop detectors to capture the high dimensional received signal features along space (over RIS geometry) and time so that more advanced DL models such as long-short term memory (LSTM) model can be leveraged.

Abbreviations

5G: Fifth-generation; 6G: Sixth-generation; PHY: Physical layer; MIMO: Multiple-input multiple-output; RIS: Reconfigurable intelligent surfaces; SISO: Single-input single-output; RAN: Radio access network; MAC: Medium access control; DL: Deep learning; RNN: Recurrent neural network; OFDM: Orthogonal frequency division multiplexing; CNN: Convolutional neural network; DRL: Deep reinforcement learning; EXP3: Exponential-weight algorithm for exploration and exploitation; FPL: Follow the perturbed leader; PL: Power law; RF: Radio frequency; LNA: Low-noise amplifier; LSTM: Long-short term memory;

Acknowledgements

This work has been funded by the Natural Sciences and Engineering Research Council of Canada, Canadian Foundation for Innovation, and the CMC Microsystems. The authors would also like to thank Philippe Massicotte for his help in the CPU/GPU setups.

Authors' contributions

MAO and DTAL conceived and designed the study, performed the experiments, and wrote the paper. All authors made suggestions for the experiments and how the data should be interpreted. All authors read, revised and approved the manuscript. MAO and DM oversaw the entire paper submission process.

Funding

This study received no external funding.

Availability of data and materials

Data sharing is not applicable to this article as no datasets were generated during the current study. All data generated or analyzed during this study are included in this article.

Declarations

Competing interests

The authors declare that they have no competing interests.

Author details

¹ M. Ahmed Ouameur and D. Massicotte are with the Université du Québec à Trois-Rivières, Department of Electrical and Computer Eng., 3351 Boul. des Forges, Trois-Rivières, Qc, Canada, G9A 5H7.

messaoud.ahmed.ouameur@uqtr.ca and daniel.massicotte@uqtr.ca

² D.T.A. Lê is with the Faculty of Information Technology, VNU-HCM University of Science, Vietnam. 20c14001@student.hcmus.edu.vn.

³ G. Jeon is with the Department of Embedded Systems Engineering, College of Information Technology, Incheon National University, Incheon, South Korea. gjeon@inu.ac.kr.

⁴ F.A.P. De Figueiredo is with the National Institute of Telecommunications, Santa Rita do Sapucaí - Minas Gerais, Brazil. felipe.figueiredo@inatel.br.

REFERENCES

- [1] Di Renzo, M. *et al.* (2020). Smart Radio Environments Empowered by Reconfigurable Intelligent Surfaces: How It Works, State of Research, and The Road Ahead. *IEEE Journal on Selected Areas in Communications*, 38(11), 2450-2525.
- [2] Basar, E. and Yildirim, I. (2020). SimRIS Channel Simulator for Reconfigurable Intelligent Surface-Empowered Communication Systems. *IEEE Latin-American Conference on Communications (LATINCOM)* (pp. 1-6).
- [3] Basar, E. (2019). Transmission through large intelligent surfaces: A new frontier in wireless communications. *Proc. Eur. Conf. Netw. Commun. (EuCNC)* (pp. 112–117).
- [4] Huang, C., Zappone, A., Alexandropoulos, G.C., Debbah, M., and Yuen, C. (2019). Reconfigurable intelligent surfaces for energy efficiency in wireless communication. *IEEE Transactions on Wireless Communications*, 18(8), 4157–4170.
- [5] Hou, T., Liu, Y., Song, Z., Sun, X., Chen, Y., and Hanzo, L. (2019). MIMO assisted networks relying on large intelligent surfaces: A stochastic geometry model. arXiv:1910.00959. [Online]. Available: <http://arxiv.org/abs/1910.00959>
- [6] Wu, Q. and Zhang, R. (2019). Intelligent reflecting surface enhanced wireless network via joint active and passive beamforming. *IEEE Transactions on Wireless Communications*, 18(11), 5394–5409.
- [7] Nokia Bell Labs (2021). Toward a 6G AI-Native Air Interface, [White paper], retrieved on July 18th 2021 from https://d1p0gxnqcu0lvz.cloudfront.net/documents/Nokia_Towards_a_6G_AI-Native_Air_Interface_Article_EN_final.pdf
- [8] Samuel, N., Diskin, T., and Wiesel, A. (2017). Deep MIMO Detection. *IEEE International workshop on signal processing advances in wireless communications* (pp. 1–5).
- [9] Mosleh, S. et al. (2018). Brain-Inspired Wireless Communications: Where Reservoir Computing Meets MIMO-OFDM. *IEEE Transactions on Neural Networks Learning Systems*, 29(10), 4694–4708.
- [10] Neumann, D., Wiese, T., and Utschick, W. (2018). Learning the MMSE Channel Estimator. *IEEE Transactions on Signal Processing*, 66(11), pp. 2905–17.
- [11] Alkhateeb, A., Alex, S., Varkey, P., Li, Y., Qu, Q. and Tujkovic, D. (2018). Deep learning coordinated beamforming for highly-mobile millimeter wave systems. *IEEE Access*, 6, 37328-37348.
- [12] Taha, A., Alrabeiah, M. and Alkhateeb, A. (2021). Enabling Large Intelligent Surfaces with Compressive Sensing and Deep Learning. *IEEE Access*, 9, 44304-44321.
- [13] Chang, H. et al. (2019). Distributive Dynamic Spectrum Access through Deep Reinforcement Learning: A Reservoir Computing-based Approach. *IEEE Internet Things Journal*, 6(2), 1938–1948.
- [14] Chinchali, S. et al. (2018). Cellular Network Traffic Scheduling with Deep Reinforcement Learning. *AAAI Conf. Artificial Intelligence*.
- [15] Shafin, R. *et al.* (2020). Self-Tuning Sectorization: Deep Reinforcement Learning Meets Broadcast Beam Optimization. *IEEE Transactions on Wireless Communications*, 19(6), 4038-4053.
- [16] Alkhateeb, A. (2019). DeepMIMO: A generic deep learning dataset for millimeter wave and massive MIMO applications. *Proceeding information theory applications workshop*, San Diego, (pp. 1-8).
- [17] Martin, C.H., Peng, T. and Mahoney, M.W. (2021). Predicting trends in the quality of state-of-the-art neural networks without access to training or testing data. *Nature Communications*, 12, 1-13.
- [18] Björnson, E. and Sanguinetti, L., (2021). Rayleigh Fading Modeling and Channel Hardening for Reconfigurable Intelligent Surfaces. *IEEE Wireless Communications Letters*, 10(4), 830-834.
- [19] Walden, R.H. (1999). Analog-to-digital converter survey and analysis. *IEEE Journal of Selected Areas Communications*, 17(4), 539-550.
- [20] Mo, J., Alkhateeb, A., Abu-Surra, S. and Heath, R.W. (2017). Hybrid architectures with few-bit ADC receivers: Achievable rates and energy-rate tradeoffs. *IEEE Transactions on Wireless Communications*, 16(4), 2274_2287.

- [21] Seldin, Y., Szepesvári, C., Auer, P. and Abbasi-Yadkori, Y. (2012). December. Evaluation and Analysis of the Performance of the EXP3 Algorithm in Stochastic Environments. *European Workshop on Reinforcement Learning*, (pp. 103–116).
- [22] Hutter, M. and Poland, J. (2005). Adaptive online prediction by following the perturbed leader. *Journal of Machine Learning Research*, 6(Apr), 639–660.
- [23] Guo, W. (2020). Explainable Artificial Intelligence for 6G: Improving Trust between Human and Machine. *IEEE Communications Magazine*, 58(6), 39-45.
- [24] WeightWatcher (2018). <https://pypi.org/project/WeightWatcher/>.
- [25] Park, S., Simeone, O. and Kang, J. (2020). Meta-Learning to Communicate: Fast End-to-End Training for Fading Channels. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, (pp. 5075-5079).
- [26] Hospedales, T., Antoniou, A., Micaelli, P. and Storkey, A. (2021). Meta-Learning in Neural Networks: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1-20.
- [27] Auer, P., Cesa-Bianchi, N., Freund, Y. and Schapire, R.E. (2002). The non-stochastic multi-armed bandit problem. *SIAM Journal on computing*, 32, 48–77.

Figure Title and Legend Section

Figure 1. The system model in which transmitter-receiver communication is aided by a RIS having an $M \times M$ interaction matrix where $M = N_y \cdot N_z$

Figure 2. The ray tracing scenario ‘O1’ from [26]. The BS#3 is selected to be the RIS while the transmitter is fixed at the position of row R850 and column 90. The receiver can be positioned at any 54300 points within the x-y grid between row R1000 and R1300. These points constitute the dataset which is split into 80% training set and 20% test set

Figure 3. The achievable rate versus the number of training data of the proposed EXP3-based scheme in comparison with the reference DL reflection beamforming [12] and the reference genie-aided method (that assumes perfect knowledge of the channel) where $\bar{M} = 4$

Figure 4. The achievable rate versus the number of active elements of the proposed EXP3-based scheme in comparison with the reference DL reflection beamforming [12] and the reference genie-aided method (that assumes perfect knowledge of the channel) where the number of training data is set to 30K

Figure 5. The energy efficiency η versus the number of active elements \bar{M} of the proposed EXP3-based scheme in comparison with the reference DL reflection beamforming [12] and the reference genie-aided method where the number of training data is set to 30K.

Figure 6. The spectral energy efficiency versus the number of training samples and the number of active elements \bar{M} for the proposed FPL-based scheme

Figure 7. The DL models used with (a) EXP3/FTPL algorithm and with (b) the reference method [12] generated using Netron

Table 1. Adversarial bandit based scheme for reflection beamforming vector perdition

Table 2. Follow the perturbed leader (FPL) scheme for reflection beamforming perdition