



**HAL**  
open science

## Reinforcement Learning for Cognitive Integrated Communication and Sensing Systems

Aya Mostafa Ibrahim Ahmed, Leila Gharsalli, Stefano Fortunati, Aydin Sezgin

► **To cite this version:**

Aya Mostafa Ibrahim Ahmed, Leila Gharsalli, Stefano Fortunati, Aydin Sezgin. Reinforcement Learning for Cognitive Integrated Communication and Sensing Systems. 26th European Microwave Week (EuMW 2023), Sep 2023, Berlin, Germany. hal-04131197

**HAL Id: hal-04131197**

**<https://hal.science/hal-04131197>**

Submitted on 16 Jun 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Reinforcement Learning for Cognitive Integrated Communication and Sensing Systems

Aya Mostafa Ahmed<sup>#1</sup>, Leila Gharsalli<sup>#2</sup>, Stefano Fortunati<sup>#3</sup> Aydin Sezgin<sup>#4</sup>

<sup>1,4</sup>DKS, Ruhr University Bochum, Germany

<sup>2,3</sup>Laboratoire des signaux et systèmes, Université Paris-Saclay, CNRS, CentraleSupélec, 91190, Gif-sur-Yvette, France

<sup>3</sup>DR2I, IPSA, 94200, Ivry-sur-Seine, France

<sup>1</sup>aya.mostafaibrahimahmad@rub.de

**Abstract**—In this paper, we propose a cognitive Massive MIMO integrated sensing and communication (ISAC) system that integrates both functionalities, enabling efficient use of the congested spectrum. To achieve this, we introduce a reinforcement learning (RL) approach that involves adaptability and that is able to optimize a joint waveform for the aforementioned system to achieve multiple objectives. We demonstrate that cognitive RL can improve state-of-the-art techniques that aims at designing the joint waveform from the ground-up achieving sensing and communication trade-off. Our results show that cognitive RL can greatly enhance sensing performance without compromising the communication performance. In contrast to previous works, we assume no prior information on the sensed scene such as the number of targets or the statistics of the disturbance.

**Keywords**—Integrated communication and sensing, 6G, Generation, Reinforcement Learning, Massive MIMO

## I. INTRODUCTION

The integration of sensing capabilities is rapidly becoming a critical component of 6G networks. Meanwhile, communication and sensing typically require different waveform specifications and metrics. The relevant literature has studied three different approaches for designing the ISAC waveform. One approach is a communication-centric design that aims to utilize traditional communication waveforms such as OFDM either directly or after performing some modifications to enable sensing [1]. However, this approach might restrict the sensing performance due to poor auto-correlation properties and its susceptibility to clutter. The second approach considered is a sensing-centric design which involves adding a communication message to an existing sensing waveform [2]. Unfortunately, such approaches might suffer from low information rates and spectral efficiency [3]. The third approach is a joint design that aims at creating a totally new waveform from the ground-up to achieve both radar and communication metrics simultaneously [4], [5]. This approach shows a great potential to achieve a balance between both systems to meet 6G demands. It must be noted however that, most of those approaches assumed prior information on the target locations in the sensing scene, which is not realistic. Other approaches as [6] addressed this problem but using a sensing-centric waveform. In [4], the authors presented a weighted optimization approach to achieve a flexible balance between the performance of radar and communications. The main objective of this optimization was to minimize

multi-user interference (MUI) while simultaneously generating a radar beam pattern towards pre-known angle locations. The authors demonstrated that their proposed design led to notable improvements in communication performance, but at the expense of some impact on radar performance. In this paper, we aim to tackle the joint waveform design, where no prior knowledge exists about the target or the surrounding environment. To achieve this, we incorporate cognitive probing of the environment, thus taking better-informed decisions. Building upon the trade-off design proposed by the authors in [4], we enhance the sensing performance using RL. Our approach could significantly enhance the detection of multiple targets within complex and unknown environments compared to the literature. Thus, we prove that by combining the cognitive probing approach with the trade-off design, we can develop a more robust and adaptable ISAC system that is better equipped to handle real-world scenarios.

## II. SYSTEM MODEL

Consider a MIMO ISAC system, which transmits communication symbols to  $K$  single antenna downlink users and simultaneously serves as monostatic cognitive radar. This joint system is capable of transmitting the same waveform to detect multiple targets and serve those users at the same time.

### A. Radar Signal Model

We assume a colocated MIMO radar, with  $N_t$  transmit and  $N_r$  receive antennas. Both the transmit and receive arrays are uniformly linear arrays (ULA), with steering vectors:  $\mathbf{a}_t(\theta)$  and  $\mathbf{a}_r(\theta)$ , respectively, where  $\theta$  is the target direction. Hence,  $\mathbf{a}_r(\theta) = [1, e^{-j\nu}, \dots, e^{-jr(N-1)\nu}]^T$  where  $\nu = \sin(\theta)$ ,  $r = \frac{2\pi d}{\lambda_c}$ ,  $d$  is the spacing between the antennas and  $\lambda_c$  is the operating frequency. Let  $\mathbf{x}_l \in \mathbb{C}^{N_t}$  be the narrow band transmitted signal, where  $l = 1, \dots, L$ , where  $L$  being the code length [7]. Thus the radar received signal at time instant  $l$  is expressed as

$$\mathbf{y}_l = \alpha \mathbf{a}_r(\theta) \mathbf{a}_t(\theta)^T \mathbf{x}_l + \mathbf{c}_l \quad (1)$$

where  $\mathbf{y}_l \in \mathbb{C}^{N_r}$ .  $\alpha \in \mathbb{C}$  is a deterministic unknown variable that accounts for the target radar cross section (RCS) and the two-way path loss,  $\mathbf{c}_l \in \mathbb{C}^{N_r}$  is the random disturbance vector. Thus, stacking the  $L$  instants in a vector form yields:

$$\mathbf{y} = \alpha \mathbf{a}_r(\theta) \mathbf{a}_t(\theta)^H \mathbf{X} + \mathbf{c} \quad (2)$$

where  $\mathbf{y}, \mathbf{c} \in \mathbb{C}^{N_r \times L}$  and  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_L] \in \mathbb{C}^{N_t \times L}$ . After correlating the received signal with the transmitted waveform and vectorizing the output, then (2) is rewritten as a function of the transmit covariance matrix  $\mathbf{R}$  as [8]

$$\mathbf{Y}_r = \alpha \mathbf{a}_r(\theta) \otimes \mathbf{a}_t(\theta)^H \mathbf{R} + \mathbf{C}_r, \quad (3)$$

where  $\mathbf{Y}_r$  and  $\mathbf{C}_r \in \mathbb{C}^N$ , where  $N = N_r N_t$ . Thus, the beam pattern produced by the transmitted waveforms in direction  $\theta$  can be expressed as  $B(\theta) = \mathbf{a}_t^H(\theta) \mathbf{R} \mathbf{a}_t(\theta)$ . It is further assumed that  $\mathbf{Y}_r$  is processed by a bank of  $B$  spatial filters, each tuned to a certain angle cell  $b$ . Afterward, hypothesis testing is performed in each angle cell. Let  $\mathbf{h} = \mathbf{a}_r(\theta) \otimes \mathbf{a}_t(\theta)^H \mathbf{R}$ , then the hypothesis testing problem can be written as

$$\begin{aligned} H_0: \quad & \mathbf{Y}_b^p = \mathbf{C}_b^p \quad p = 1, \dots, P \\ H_1: \quad & \mathbf{Y}_b^p = \alpha_b^p \mathbf{h}_{p,b} + \mathbf{C}_b^p \quad p = 1, \dots, P. \end{aligned} \quad (4)$$

The decision statistic is chosen as robust Wald-type detection which ensures CFAR property as  $N \rightarrow \infty$  [9] as

$$\Lambda_{p,b} = \frac{2|\mathbf{h}_{p,b}^H \mathbf{Y}_b^p|^2}{\mathbf{h}_{p,b}^H \hat{\mathbf{\Gamma}}_{p,b} \mathbf{h}_{p,b}} \quad (5)$$

where  $\hat{\mathbf{\Gamma}}_{p,b}$  is the estimated covariance matrix of  $\mathbf{C}_b^p$ , defined in [10, eq.(38)]. To distinguish between  $H_0$  and  $H_1$ , the detector evaluates (5) against a threshold  $\lambda = -2 \ln(P_{FA})$  in each angle bin, where  $P_{FA} = 10^{-4}$  is the false alarm rate.

### B. Communication Signal Model

The downlink received signal at the users can be expressed as [4]:

$$\mathbf{Y}_c = \mathbf{H}\mathbf{X} + \mathbf{N}, \quad (6)$$

where  $\mathbf{H} \in \mathbb{C}^{K \times N_t}$ ,  $\mathbf{N} \in \mathbb{C}^{K \times L}$ . Here  $\mathbf{X}$  is the dual waveform transmitted for both radar and communication purposes, thus while  $L$  previously denoted the number of snapshots within a radar pulse  $p$ , in (6),  $L$  denotes the number of symbols in the communication frame. Here we assume that channel  $\mathbf{H}$  is a flat fading Rayleigh channel that remains unchanged within the communication frame or radar pulse. Therefore, if the desired constellation matrix for the  $K$  users is given as  $\mathbf{S} \in \mathbb{C}^{K \times L}$ . Then (6) can be rewritten as

$$\mathbf{Y}_c = \mathbf{S} + \underbrace{(\mathbf{H}\mathbf{X} - \mathbf{S})}_{\text{MUI}} + \mathbf{N}, \quad (7)$$

where the underlined term denotes the MUI. It is worth noting that minimizing the MUI power is directly linked to increasing the achievable sum-rate for all users [4].

## III. JOINT COMMUNICATION AND RADAR WAVEFORM DESIGN

The metrics used to evaluate the communication and sensing performance are distinct due to the nature of both systems. The purpose of communication is to convey as much information as possible to a certain receiver accurately, while the radar objective is to extract specific information from the

returned echoes. Unlike communication, the useful information here is in the received signal, not the transmitted waveform. From the previously discussed system models of both systems, it can be depicted that one way to design the waveform is to formulate an optimization problem that aims at minimizing the MUI power, while at the same time ensuring a specific beam pattern for detecting the radar targets. Therefore, we consider the trade-off design model proposed in [4] to optimize the transmitted waveform to fulfill both objectives as follows:

$$\begin{aligned} \mathcal{P}_0: \quad & \min_{\mathbf{X}} \rho \|\mathbf{H}\mathbf{X} - \mathbf{S}\|_F^2 + (1 - \rho) \|\mathbf{X} - \mathbf{X}_0\|_F^2 \\ \text{s.t.} \quad & \frac{1}{L} \|\mathbf{X}\|_F^2 = P_T, \end{aligned} \quad (8)$$

where  $0 \leq \rho \leq 1$  is a weighing factor that allows one to favor one system over the others in the ISAC framework. While the first term in  $\mathcal{P}_0$  minimizes the MUI, the second term aims at finding  $\mathbf{X}$  that matches  $\mathbf{X}_0$ , which can be obtained by solving the following optimization problem

$$\begin{aligned} \mathcal{P}_1: \quad & \min_{\mathbf{X}_0} \|\mathbf{H}\mathbf{X}_0 - \mathbf{S}\|_F^2 \\ \text{s.t.} \quad & \frac{1}{L} \mathbf{X}_0 \mathbf{X}_0^H = \mathbf{R}_d. \end{aligned} \quad (9)$$

The optimization problem  $\mathcal{P}_1$  solves for  $\mathbf{X}_0$  through minimizing the MUI while achieving a certain desired covariance matrix  $\mathbf{R}_d$  that corresponds to a certain desired beam pattern  $B(\theta)$ . The authors in [4] proposed a solution to  $\mathcal{P}_1$ , however they assumed that the covariance matrix  $\mathbf{R}_d$  is given. Such an assumption is not practical from the radar perspective, since in a sensing scenario, the number of targets and their possible directions is usually unknown. In addition, they used a generalized likelihood ratio test (GLRT) to detect the targets within independent and identically distributed (i.i.d) Gaussian noise. This assumption might not be valid due to the dynamic and complex nature of the disturbance in radar environments. Thus, we propose an RL-based detection algorithm that cognitively probes the environment using the ISAC waveform in  $\mathcal{P}_0$  under the following assumptions 1) no prior information is assumed on the targets locations or number 2) no statistical characterization of the clutter is assumed.

### A. RL-based Cognitive Waveform design

RL is a machine learning technique that facilitates an agent's ability to achieve a designated goal by interacting with and learning from the surrounding environment through trial and error. In our case, the ISAC system is the agent, which receives continuous feedback from the environment. Based on those, the agent determines its next *action*  $a_p$ . Here, we use SARSA algorithm, where the agent learns to select its next *action*  $a_p$  based on the current *state* of the environment  $s_p$ , then receives a reward  $r_p$  after each action taken. The algorithm keeps updating the estimate of the optimal policy based on the *state-action-reward-state-action* (SARSA) transitions during its interaction with the environment. Such policy is determined based on the  $Q$ -value of the current state-action transition. The  $Q$ -value represents the expected, discounted, and cumulative

reward that the agent receives by taking action  $a_p$ , in particular state  $s_p$  following certain policy  $\pi$ , which is defined as

$$Q(s_p, a_p) \leftarrow Q(s_p, a_p) + \alpha (r_{p+1} + \gamma Q(s_{p+1}, a_{p+1}) - Q(s_p, a_p)). \quad (10)$$

The learning rate  $\alpha \in [0, 1]$  determines the extent to which new experiences should overwrite old ones. The discounted factor  $\gamma$  on the other hand controls the influence of future rewards in the decision-making process. Hence, SARSA is used to navigate through the unknown environment of ISAC system, while updating the state-action matrix  $\mathbf{Q} \in \mathbb{R}^{(M+1) \times (M+1)}$ , where  $M$  is the maximum number of detectable targets. Consequently, the joint waveform is continuously adapted based on (10). Next, the SARSA parameters are further explained.

#### 1) The set of states

A state  $s_p$  represents the current situation of the unknown environment. Here, the state is defined in terms of the detection statistic for each angle cell  $b$ , defined in (5), thus, we define the state after having sent a pulse at the time  $p$  as the number of detections using the Wald test such that

$$s_p = \sum_{b=1}^B \bar{\Lambda}_b^p, \quad (11)$$

where

$$\bar{\Lambda}_b^p = \begin{cases} 1 & \Lambda_b^p > \lambda_\Lambda \\ 0 & \text{otherwise.} \end{cases} \quad (12)$$

Thus, the state space can be defined as  $\mathcal{S} \triangleq \{0, \dots, M\}$  [10].

#### 2) The set of actions

The action is defined as transmitting joint ISAC waveform based on certain communication channel  $H$  and  $\mathbf{R}_d$ . Without the loss of generality, we assume perfectly known  $\mathbf{H}$ , thus we focus on designing  $\mathbf{R}_d$  in  $\mathcal{P}_1$  according to the following problem [11]:

$$\begin{aligned} \mathcal{P}_2 : \max_{\mathbf{R}_d} \quad & \text{tr}(\mathbf{R}_d \hat{\mathbf{B}}) \\ \text{subject to} \quad & \text{tr}(\mathbf{R}_d) = P_T \\ & \mathbf{R}_d \geq 0. \end{aligned} \quad (13)$$

where  $\hat{\mathbf{B}} = \sum_{b=1}^{\hat{b}} \mathbf{a}(\theta_b) \mathbf{a}^*(\theta_b)$ . Thus,  $\mathcal{P}_2$  namely maximizes the total power at the estimated angle bins's locations. Hence, the role of SARSA is to define  $\Theta_{\hat{b}} = \{\hat{\theta}_1, \dots, \hat{\theta}_{\hat{b}}\}$  and  $\hat{\theta}$  is the estimated angle bin of the target. Furthermore,  $\Theta_p$  is calculated based on the highest  $\hat{b}$  values of  $\Lambda_{p,b}$  in (5). Accordingly, we can define the action as  $a_p \in \mathcal{A} = \{\Theta_i | i \in \{1, 2, \dots, M\}\}$ . After finding  $\Theta_{\hat{b}}$ ,  $\mathcal{P}_2$  will have a closed form solution defined in [11], eq. (14). In order to evaluate how good those actions are, a reward is given to the agent. Thus the agent's goal is continuously maximize the cumulative reward [12]. It is defined as

$$r_{p+1} = \sum_{b=1}^{s_p} \hat{P}_{D_b}^p - \sum_{j=1}^{B-s_k} \hat{P}_{D_j}^p, \quad (14)$$

where  $\hat{P}_{D_b}^p$  is the estimated probability of detection, which can be defined in closed form expression as  $N \rightarrow \infty$  in [9]. The policy  $\pi_p(s_p)$  used to maximize the reward is quasi  $\epsilon$  greedy [13] where the algorithm chooses  $a_p^{(\text{greedy})} \triangleq \arg \max_{a \in \mathcal{A}} \mathbf{Q}(s_{p+1}, a)$  with a probability (w.p) of  $1 - \epsilon$ . Meanwhile, another random action  $a_{\text{rnd}}$  (excluding  $a_p^{(\text{greedy})}$ ) is chosen from the set <sup>1</sup>  $\mathcal{A}'(s_p^j) \triangleq \{\Theta_i, i = j, \dots, M\}$  w.p  $\epsilon$ . Such definition of the policy would allow the exploration of random actions based on  $\epsilon$ .

---

#### Algorithm 1 SARSA

---

Input:  $\mathbf{H}$ ,  $\mathbf{S}$ ,  $\rho$ , and  $P_T$   
Initialize  $s_0 = 1$ ,  $a_0 = 1$ ,  $P = 50$ ,  $\mathbf{Q} = \mathbf{0}$  and  $\mathbf{X}_p = \mathbf{I}$   
**repeat** for each pulse  $p$ :  
    Take action  $a_p$  by transmitting waveform  $\mathbf{X}_p$   
    Acquire the received signal  $\mathbf{Y}_p^b, \forall b$   
    Solve for the decision statistic in (5)  
    Calculate  $s_{p+1}$  from (11) and  $r_{p+1}$  as in (14)  
    Choose action  $a_{p+1}$  with quasi  $\epsilon$  greedy, identify  $\Theta_{\hat{b}}$   
    Solve for  $\mathbf{R}_d$  in  $\mathcal{P}_2$   
    Solve  $\mathbf{X}_0$  in  $\mathcal{P}_1$ ,  $\mathbf{X}$  in  $\mathcal{P}_0$  as [4] then transmit  
    Update  $Q(s_p, a_p)$  as in (10),  $s_k \leftarrow s_{p+1}; a_p \leftarrow a_{p+1}$   
**until** Observation time ends

---

## IV. RESULTS

In our simulations, we consider a massive MIMO base station (BS), which communicates messages of length  $L = 10$  to  $K = 4$  users whose channel entries  $H_{i,j} \in \mathcal{CN}(0, 1)$ . Similar to [4],  $\mathbf{H}$  is assumed to be perfectly estimated. Furthermore, the constellation is chosen to be QPSK alphabet of unit power. The BS simultaneously tries to detect 4 targets using  $P = 50$  pulses. The sensing scene is divided into total of  $B = 20$  angle bins, where the angle grid is defined as  $\nu = [-0.5 : 0.45]$ . The four targets are at locations  $\nu = \{-0.2, 0, 0.2, 0.3\} \subset \nu$ , with SNR =  $[-30, -20, -10, -20]$ dB respectively. We further assumed that the targets are masked within unknown disturbance modeled as AR (6) with t-distributed i.i.d innovations defined as in [10]. In Fig 1, the probability of detection is averaged across the 50 pulses and calculated across the number of spatial channels  $N$  for the targets with the least SNR at  $\nu = -0.2$  and  $\nu = 0$ . Here,  $\rho = 0.8$ , which according to (8) gives the main weight to the communication part of the design. We compare our algorithm to the directional and omnidirectional trade-off approaches proposed in [4]. To ensure a fair comparison, we adapted those approaches to a cognitive design, where the BS adaptively modifies  $\mathbf{X}$  based on the same detection statistic in (5) without the RL component. Despite this adaptation, our RL algorithm still demonstrated superior performance for the target at  $\nu = 0.2$ , while the other algorithms failed to detect it altogether, even as  $N$  increased asymptotically.

<sup>1</sup>Please note that  $s^i$  with superscript denotes the value of the state, while  $s_p$  with subscript denotes the state at pulse  $p$  (i.e.,  $s_0 = s^0$  means the state at  $p = 0$  has the value of  $s^0$ )

## V. CONCLUSION

In this work, we proposed a cognitive waveform design for ISAC system, that is capable of serving communication users while detecting multiple targets in an unknown environment. The approach improves the trade-off design proposed in the literature, by employing cognitive RL, to achieve the balance between sensing and communications. We prove that our approach significantly enhances the system's sensing performance, where it could detect fading targets of very low SNR, without trading the communication performance.

## ACKNOWLEDGMENT

This work is funded by the German Federal Ministry of Education and Research (BMBF) in the course of the 6GEM Research Hub under Grant 16KISK037.

## REFERENCES

- [1] Y. Liu, G. Liao, Z. Yang, and J. Xu, "Design of integrated radar and communication system based on mimo-ofdm waveform," *Journal of Systems Engineering and Electronics*, vol. 28, no. 4, pp. 669–680, 2017.
- [2] T. Huang, X. Xu, Y. Liu, N. Shlezinger, and Y. C. Eldar, "A dual-function radar communication system using index modulation," in *2019 IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 2019, pp. 1–5.
- [3] W. Zhou, R. Zhang, G. Chen, and W. Wu, "Integrated sensing and communication waveform design: A survey," *IEEE Open Journal of the Communications Society*, vol. 3, pp. 1930–1949, 2022.
- [4] F. Liu, L. Zhou, C. Masouros, A. Li, W. Luo, and A. Petropulu, "Toward dual-functional radar-communication systems: Optimal waveform design," *IEEE Transactions on Signal Processing*, vol. 66, no. 16, pp. 4264–4279, 2018.
- [5] X. Liu, T. Huang, N. Shlezinger, Y. Liu, J. Zhou, and Y. C. Eldar, "Joint transmit beamforming for multiuser mimo communications and mimo radar," *IEEE Transactions on Signal Processing*, vol. 68, pp. 3929–3944, 2020.
- [6] W. Zhai, X. Wang, X. Cao, M. S. Greco, and F. Gini, "Reinforcement learning based dual-functional massive mimo systems for multi-target detection and communications," *IEEE Transactions on Signal Processing*, pp. 1–15, 2023.
- [7] W. Wu, B. Tang, and X. Wang, "Constant-modulus waveform design for dual-function radar-communication systems in the presence of clutter," *IEEE Transactions on Aerospace and Electronic Systems*, pp. 1–14, 2023.
- [8] S. Ahmed and M.-S. Alouini, "Mimo-radar waveform covariance matrix for high sinr and low side-lobe levels," *IEEE Transactions on Signal Processing*, vol. 62, no. 8, pp. 2056–2065, 2014.
- [9] S. Fortunati, L. Sanguinetti, F. Gini, M. S. Greco, and B. Himed, "Massive mimo radar for target detection," *IEEE Transactions on Signal Processing*, vol. 68, pp. 859–871, 2020.
- [10] A. M. Ahmed, A. A. Ahmad, S. Fortunati, A. Sezgin, M. S. Greco, and F. Gini, "A reinforcement learning based approach for multitarget detection in massive mimo radar," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 57, no. 5, pp. 2622–2636, 2021.
- [11] P. Stoica, J. Li, and Y. Xie, "On probing signal design for mimo radar," *IEEE Transactions on Signal Processing*, vol. 55, no. 8, pp. 4151–4161, 2007.
- [12] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. The MIT Press, 2018. [Online]. Available: <http://incompleteideas.net/book/the-book-2nd.html>
- [13] F. Lisi, S. Fortunati, M. S. Greco, and F. Gini, "Enhancement of a state-of-the-art rl-based detection algorithm for massive mimo radars," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 58, no. 6, pp. 5925–5931, 2022.

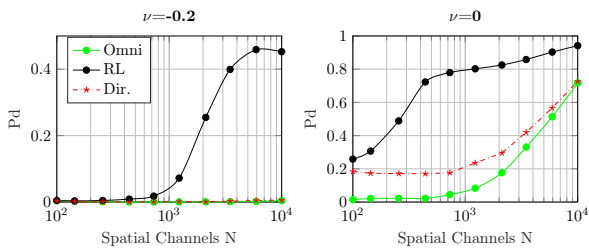


Figure 1. Pd vs Spatial channels,  $N = N_t N_r$ .

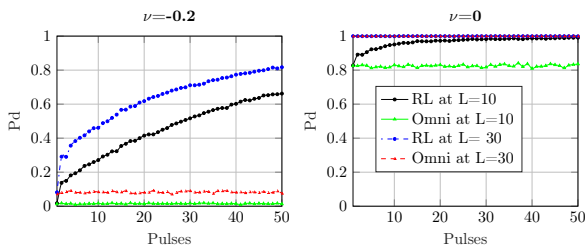


Figure 2. Pd vs Pulses.

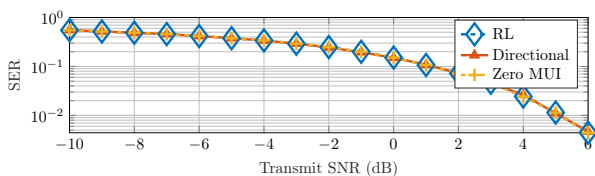


Figure 3. SER vs SNR.

Notably, both approaches perform better for  $\nu = 0$  due to its higher SNR, reaching the same performance at  $N=10^4$ , however our RL approach still outperformed both approaches in this scenario as well. We further explore this performance in Fig. 2 using  $N_t = N_r = 100$ , we compare the  $P_d$  as function of the number of pulses. We only chose to compare against the omnidirectional approach, since it performs similar to the directional one symptomatically. It can be noticed for the first target, the  $P_d$  increases over time, due to the nature of RL which learns from the environment using trial and error. In this scenario, the effect of increasing  $L$  is explored. Similar observations can be drawn for both targets, however it can be noticed that as  $L$  increases to 30 the  $P_d$  also increases, leading to full detection for the target with the lower SNR, since more snapshots are considered in this case.

The proposed algorithm focuses on improving the sensing performance only while keeping the communication aspect unchanged. We further simulated the symbol error rate (SER) of the RL and directional approaches at  $N_t = 100$ , it can be noticed from Fig.3 that both approaches yield the same performance, which aligns with the zero MUI case as well, since the communication component in (8) is weighted by 0.8. Notably, RL only improved the detection performance leaving room for further enhancements in communication performance as future work.