

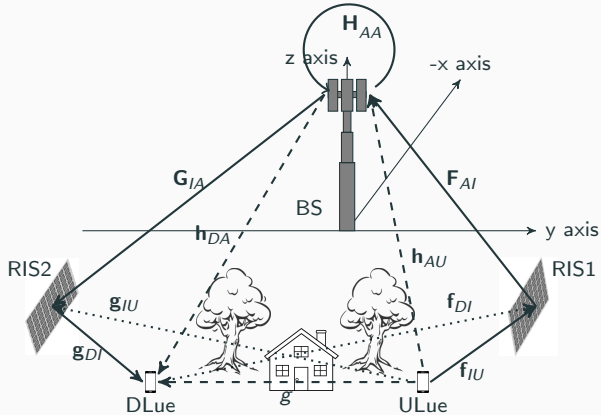
A DRL Approach for RIS-Assisted Full-Duplex UL and DL Transmission: Beamforming, Phase Shift and Power Optimization

Speaker: Dr. Nancy Nayak

Imperial College London

- Machine learning extracts valuable insights from data
- ML methods learn from data and minimize the reliance on classical estimation techniques
- Existing methods use CSI and the knowledge of self-interference
- Proposed: a **two-stage DRL framework** to solve RIS phase-shift, and beamformers and transmit powers
- Strengths: learns directly from data without CSI, enables end-to-end learning, handles complex non-linear environments, solves non-convex problems in one step, and flexibly adapts to changing optimization objectives.
- Reduces signaling overhead

RIS assisted Full Duplex Communication¹



¹Nayak, Nancy, Sheetal Kalyani, and Himel A. Suraweera. "A DRL approach for RIS-assisted full-duplex UL and DL transmission: Beamforming, phase shift and power optimization." IEEE Transactions on Wireless Communications (2024).

System model

- Full-duplex (FD) setup, one FD BS, one Half duplex (HD) ULue, one HD DLue, and two Reconfigurable Intelligent Surfaces (RIS) to facilitate communication when the users are not in LoS with the BS
- BS has a uniform linear antenna (ULA) array with M_t transmit antenna elements and M_r receive antenna elements
- ULue and DLue are single-antenna HD user
- Both the RISs are deployed as Uniform Planar Antenna (UPA)
- RIS1 and RIS2 has $N_1 = N_{1h}N_{1v}$ and $N_2 = N_{2h}N_{2v}$ reflecting elements
- The direct paths are blocked therefore has high pathloss
- User to user interference and inter-RIS interference is present, but paths have high pathloss
- DLue is HD, therefore it cannot transmit and receive at the same time

- The diagonal phase-shift matrices Θ_U and Θ_D of the two RISs are

$$\begin{aligned}\Theta_U &= \text{diag}\{\bar{\Theta}_U\} = \text{diag}\{\phi_{U1} \dots \phi_{UN_1}\}, \text{ and} \\ \Theta_D &= \text{diag}\{\bar{\Theta}_D\} = \text{diag}\{\phi_{D1} \dots \phi_{DN_2}\}\end{aligned}\tag{1}$$

with $\phi_n = e^{j\theta_n}$

- s_U denotes the transmit signal from the ULue
- $p_U > 0$ denotes the transmit power of the ULue
- s_D denotes the transmit signal of the BS
- $p_A > 0$ denotes the transmit power of the BS

System model

- The received signal at the BS is given by

$$y_A = \mathbf{w}_R \mathbf{h}_{AU} \sqrt{p_U} s_U + \mathbf{w}_R \mathbf{F}_{AI} \mathbf{\Theta}_U \mathbf{f}_{IU} \sqrt{p_U} s_U + \mathbf{w}_R \mathbf{G}_{IA}^T \mathbf{\Theta}_D \mathbf{g}_{IU} \sqrt{p_U} s_U \\ + \underbrace{\mathbf{w}_R \mathbf{H}_{AA} \mathbf{w}_T \sqrt{p_A} s_D}_{\text{residual SI}} + \mathbf{w}_R n_A, \quad (2)$$

- The signal received by the DLue is

$$y_D = \mathbf{h}_{DA} \mathbf{w}_T \sqrt{p_A} s_D + \mathbf{g}_{DI} \mathbf{\Theta}_D \mathbf{G}_{IA} \mathbf{w}_T \sqrt{p_A} s_D + \mathbf{f}_{DI} \mathbf{\Theta}_U \mathbf{F}_{AI}^T \mathbf{w}_T \sqrt{p_A} s_D \\ + \underbrace{\mathbf{g}_{DI} \mathbf{\Theta}_D \mathbf{g}_{IU} \sqrt{p_U} s_U + \mathbf{f}_{DI} \mathbf{\Theta}_U \mathbf{f}_{IU} \sqrt{p_U} s_U}_{\text{inter-RIS}} + \underbrace{g \sqrt{p_U} s_U}_{\text{inter-user}} + n_D. \quad (3)$$

- Here transmit and receiver beamformers are denoted by \mathbf{w}_T and \mathbf{w}_R respectively
- Highlighted term leads to a very high level of interference if not canceled

- The SINR at the BS, γ_{BS} and the DL user, γ_{DL} are given by

$$\begin{aligned}\gamma_{BS} &= \frac{p_U \|\mathbf{w}_R(\mathbf{h}_{AU} + \mathbf{F}_{AI} \mathbf{\Theta}_U \mathbf{f}_{IU} + \mathbf{G}_{IA}^T \mathbf{\Theta}_D \mathbf{g}_{IU})\|^2}{p_A \|\mathbf{w}_R \mathbf{H}_{AA} \mathbf{w}_T\|^2 + \mathbf{w}_R^2 \sigma_A^2}, \text{ and} \\ \gamma_{DL} &= \frac{p_A \|\mathbf{h}_{DA} + \mathbf{g}_{DI} \mathbf{\Theta}_D \mathbf{G}_{IA} \mathbf{w}_T + \mathbf{f}_{DI} \mathbf{\Theta}_U \mathbf{F}_{AI}^T \mathbf{w}_T\|^2}{p_U \|(g + \mathbf{g}_{DI} \mathbf{\Theta}_D \mathbf{g}_{IU} + \mathbf{f}_{DI} \mathbf{\Theta}_U \mathbf{f}_{IU})\|^2 + \sigma_D^2},\end{aligned}\tag{4}$$

respectively.

- Accordingly, the data rate at the DLue and the BS are given by

$$\begin{aligned}r_{DL} &= \log_2(1 + \gamma_{DL}), \text{ and} \\ r_{BS} &= \log_2(1 + \gamma_{BS}).\end{aligned}\tag{5}$$

Objective

The optimization problem is formulated as follows:

$$\begin{aligned} \mathcal{P}_1 : \quad & \max_{\Theta_D, \Theta_U, \mathbf{w}_T, \mathbf{w}_R, p_A, p_U} r_{BS} + r_{DL} \\ & \text{s.t. } p_A^{max} \geq p_A \geq 0, \quad p_U^{max} \geq p_U \geq 0, \\ & |\phi_n| = 1, 1 \leq n \leq N_1, 1 \leq n \leq N_2 \end{aligned} \tag{6}$$

- Such optimization problems are typically relaxed and broken into smaller subproblems
- Either assumes negligible residual self-interference due to excellent SI mitigation technique
- Or assumes the presence of very little residual SI and then
 - cancels the residual SI by using beamformers
 - however, designing beamformers needs CSI
- What happens if the CSI is noisy and the residual SI is high/unknown?
- Proposed solution is **two-stage Deep Reinforcement Learning (DRL) algorithm** that does not need any CSI

DRL based two-stage algorithm

- Main **challenge in the FD communication** system: the **self-interference (SI)** imposed by the transmit antenna on the receive antenna of the BS
- If SI mitigation scheme is not good, the residual SI has high power
- As the **DRL agent learns from feedback**, it is difficult for the DRL agent to learn anything **if the received signal has too much interference** due to the high residual SI
- So first stage is to **cancel a major part of the SI interference by sending a pilot symbol** (next slide)
- Second stage is to **feed the data to an Intelligent agent** which learns to predict the RIS phaseshifts, beamformers and the transmit powers (agent: DRL algorithm)

First stage: Least square-based SI-cancellation (LSSIC)

- The signal due to SI can be canceled with an estimate of $\mathbf{w}_R \mathbf{H}_{AA} \mathbf{w}_T$ denoted by \hat{h} as

$$y_A = \mathbf{w}_R \mathbf{h}_{AU} \sqrt{p_U} s_U + \mathbf{w}_R \mathbf{F}_{AI} \mathbf{\Theta}_U \mathbf{f}_{IU} \sqrt{p_U} s_U + \mathbf{w}_R \mathbf{G}_{IA}^T \mathbf{\Theta}_D \mathbf{g}_{IU} \sqrt{p_U} s_U \\ + \underbrace{(\mathbf{w}_R \mathbf{H}_{AA} \mathbf{w}_T - \hat{h}) \sqrt{p_A} s_D}_{\text{residual SI}} + \mathbf{w}_R n_A. \quad (7)$$

- At every epoch, the scalar \hat{h} is estimated by sending a pilot signal $s_D^p \in \mathcal{C}$ at the BS from the transmitter to the receiver antenna
- The corresponding received signal $y_A^p \in \mathcal{C}$ at the receiver antenna of BS can be expressed as

$$y_A^p = \mathbf{w}_R \mathbf{H}_{AA} \mathbf{w}_T \sqrt{p_A} s_D^p + v_A, \quad (8)$$

where $v_A \in \mathcal{C}$ is the AWGN and the scalar $\mathbf{w}_R \mathbf{H}_{AA} \mathbf{w}_T$ needs to be estimated.

Contd. First stage: LSSIC

- To estimate $\mathbf{w}_R \mathbf{H}_{AA} \mathbf{w}_T$, we need to minimise the error $J(h)$ where

$$\begin{aligned} J(h) &= (\overline{y_A^p} - \sqrt{p_A} \overline{s_D^p} h)(y_A^p - \sqrt{p_A} s_D^p h), \\ &= \overline{y_A^p} y_A^p - 2\sqrt{p_A} h \overline{y_A^p} s_D^p + p_A h^2 \overline{s_D^p} s_D^p. \end{aligned} \tag{9}$$

where \bar{x} denotes the conjugate transpose of x .

- By taking the derivative of $J(h)$ with respect to h and equating it to zero to obtain \hat{h} ,

$$\begin{aligned} \frac{\partial J(h)}{\partial h} &= 0 - 2\sqrt{p_A} \overline{y_A^p} s_D^p + 2p_A h \overline{s_D^p} s_D^p = 0, \\ \text{therefore, } \hat{h} &= \frac{1}{\sqrt{p_A}} (\overline{s_D^p} s_D^p)^{-1} \overline{s_D^p} y_A^p. \end{aligned} \tag{10}$$

- The derived \hat{h} can be used in (7) to cancel a significant amount of SI.

Second stage: DRL based method

- The DRL agent located at the BS initializes the phaseshift, beamformers and transmit powers randomly, together called as **actions**
- Using these, the BS and ULue transmit and DLue and BS receives signals
- After receiving the signal at the BS, a significant amount of residual SI is cancelled² using a pilot signal (LSSIC); then the SINR at BS is obtained; used as **observation or state**
- The SINR at DLue (a scalar) is also calculated and the SINR is fed back to the agent to be used in second stage as **observation or state**
- Note, if an estimate of \mathbf{H}_{AA} is already available at the BS, \hat{h} can be estimated using this $\tilde{\mathbf{H}}_{AA}$; named as $\tilde{\mathbf{H}}_{AA}$ based SI cancellation (HSIC)

²DLue is HD, so no SI for DLue

Second stage: DRL based method

- Weighted sum of the corresponding data rates of these two SINRs are used as **reward**
- Now based on this feedback, the DRL agent predicts a different set of actions, according to which again the signals are transmitted from BS and ULue, achieving a new pair of UL and DL SINR (after LSSIC/HSIC at the first stage for UL) which takes the system to a **next state**
- The quality of the beamformers learned by the DRL agent in the second stage depends on the SI-cancelled signal from the first stage
- At the same time, with accurate learning of the beamformers, the SI-cancellation is also better

Formulating MDP

- An MDP has a **state** space \mathcal{S} , an **action** space \mathcal{A} , an initial distribution of space $p(\mathbf{s}^{\{1\}})$ and a stationary distribution for state transition that obeys Markov property i.e., $p(\mathbf{s}^{\{t+1\}}|\mathbf{s}^{\{t\}}, \mathbf{a}^{\{t\}}) = p(\mathbf{s}^{\{t+1\}}|\mathbf{s}^{\{t\}}, \mathbf{a}^{\{t\}}, \dots, \mathbf{s}^{\{1\}}, \mathbf{a}^{\{1\}})$ and a **reward** function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$.
- The algorithm (deployed at BS) gives action $\mathbf{a}^{\{t\}}$ based on the state $\mathbf{s}^{\{t\}}$ generated at a previous time step
- The environment reacts to these **actions** and gives back the SINRs as the **observations/states** indicate how good the actions are
- Finally, the **reward** is calculated and fed as input to the learning agent
- The SINR observations, along with the actions at time step t give the state for time step $(t + 1)$.

MDP for our system model

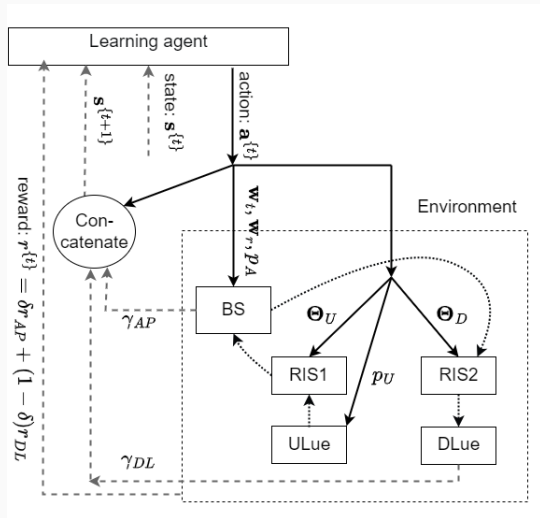


Figure 1: MDP formulation for RIS-based FD communication.

Contd. Second stage: DRL based method

- The job of the RL agent is to learn the policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ from the observations corresponding to each of the actions while maximizing the return
$$\mathbb{r}^{\{t\}}(\gamma) = \sum_{t'=t}^{\infty} \gamma^{t'-t} r^{\{t\}}(\mathbf{a}^{\{t\}}, \mathbf{s}^{\{t\}})$$
 where $\gamma \in [0, 1]$ is the discounting factor.
- The learning agent calculates the quality of action by using Q -function given by $Q^{\pi}(s, a) = \mathbb{E}[\mathbb{r}_1(\gamma) | S_1 = s, A_1 = a; \pi]$ indicating how rewarding each action a is when taken from a state s . At each timestep, the agent takes action which maximizes the Q -value
- A better agent is an algorithm that approximates the Q -value well and predicts a good action/policy
- The policy can be approximated by deep neural networks - **actor-critic method namely deep deterministic policy gradient (DDPG)**

DDPG has four neural networks:

- An **actor-network** \mathbb{A} parameterized by ω^a which predicts the action $\mathbf{a}^{\{t\}}$ based on the current state $\mathbf{s}^{\{t\}}$
- A **critic network** \mathbb{C} parameterized by ω^c which computes $Q(\mathbf{s}^{\{t\}}, \mathbf{a}^{\{t\}})$ that is essentially the quality of the action taken by actor-network \mathbb{A}
- A target actor and a target critic networks for stable updates
- The agent encourages the actor-network to take **better actions through its feedback**
- The critic network \mathbb{C} trains itself for **better prediction by observing the rewards after each action** ³
- Critic network is a feed forward network

³For more details regarding how these networks are trained, please refer to the paper.

Proposed actor network

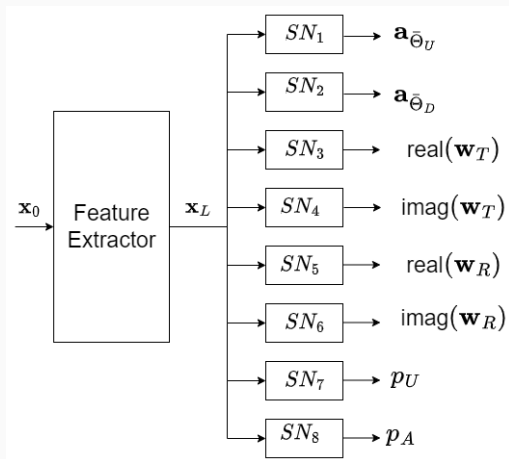


Figure 2: The proposed action predictor network. SN represents sub-network

Actor network

- The state for next time step $\mathbf{s}^{\{t+1\}}$ takes the SINR at BS and DLue $\gamma_{BS}^{\{t\}}$ and $\gamma_{DL}^{\{t\}}$, RIS phases $\bar{\Theta}_U^{\{t\}}$, $\bar{\Theta}_D^{\{t\}}$, transmit and receive beamformers $\mathbf{w}_T^{\{t\}}$ and $\mathbf{w}_R^{\{t\}}$ and transmit powers of the BS and ULue given by $p_A^{\{t\}}$ and $p_U^{\{t\}}$
- Input to actor network \mathbf{x}_0 is the state

$$\begin{aligned} \mathbf{s}^{\{t\}} = & [\gamma_{BS}^{\{t-1\}}, \gamma_{DL}^{\{t-1\}}, \theta_{11}^{\{t-1\}}, \theta_{12}^{\{t-1\}}, \dots, \theta_{1N_1}^{\{t-1\}}, \theta_{21}^{\{t-1\}}, \theta_{22}^{\{t-1\}}, \dots, \theta_{2N_2}^{\{t-1\}}, \\ & \text{real}(w_{t1}^{\{t-1\}}, w_{t2}^{\{t-1\}}, \dots, w_{tM_t}^{\{t-1\}}), \text{imag}(w_{t1}^{\{t-1\}}, w_{t2}^{\{t-1\}}, \dots, w_{tM_t}^{\{t-1\}}), \\ & \text{real}(w_{r1}^{\{t-1\}}, w_{r2}^{\{t-1\}}, \dots, w_{rM_r}^{\{t-1\}}), \text{imag}(w_{r1}^{\{t-1\}}, w_{r2}^{\{t-1\}}, \dots, w_{rM_r}^{\{t-1\}}), p_A^{\{t-1\}}, p_U^{\{t-1\}}]. \end{aligned} \quad (11)$$

- The predicted action by the actor network for time step t is given by

$$\begin{aligned} \mathbf{a}^{\{t\}} = & [\theta_{11}^{\{t\}}, \theta_{12}^{\{t\}}, \dots, \theta_{1N_1}^{\{t\}}, \theta_{21}^{\{t\}}, \theta_{22}^{\{t\}}, \dots, \theta_{2N_2}^{\{t\}}, \text{real}(w_{T1}^{\{t\}}, w_{T2}^{\{t\}}, \dots, w_{TM_t}^{\{t\}}), \\ & \text{imag}(w_{T1}^{\{t\}}, w_{T2}^{\{t\}}, \dots, w_{TM_t}^{\{t\}}), \text{real}(w_{R1}^{\{t\}}, w_{R2}^{\{t\}}, \dots, w_{RM_r}^{\{t\}}), \\ & \text{imag}(w_{R1}^{\{t\}}, w_{R2}^{\{t\}}, \dots, w_{RM_r}^{\{t\}}), p_A^{\{t\}}, p_U^{\{t\}}]. \end{aligned} \quad (12)$$

Predicted actions

- **RIS phases:** The feature \mathbf{x}_L is passed via the first two subnetworks, and the outputs are:

$$\mathbf{a}_{\bar{\Theta}_U} = \tanh(\mathbf{W}_{\bar{\Theta}_U} \mathbf{x}_L + \mathbf{b}_{\bar{\Theta}_U}), \text{ and } \mathbf{a}_{\bar{\Theta}_D} = \tanh(\mathbf{W}_{\bar{\Theta}_D} \mathbf{x}_L + \mathbf{b}_{\bar{\Theta}_D}). \quad (13)$$

The \tanh is used to get the normalized actions between $[-1, +1]$ are then shifted and scaled to take values in $[0, 2\pi]$.

- **Beamformers:** The in-phase and quadrature part of beamforming vectors can take any value between $[-1, +1]$, so the action corresponding to the in-phase and quadrature components of the transmit beamforming vector ⁴

$$\mathbf{a}_{M_t, I} = \text{real}(\mathbf{w}_T) = \tanh(\mathbf{W}_{2, M_t, I} \text{ReLU}(\mathbf{W}_{1, M_t, I} \mathbf{x}_L + \mathbf{b}_{1, M_t, I}) + \mathbf{b}_{2, M_t, I}) \text{ and}$$
$$\mathbf{a}_{M_t, Q} = \text{imag}(\mathbf{w}_T) = \tanh(\mathbf{W}_{2, M_t, Q} \text{ReLU}(\mathbf{W}_{1, M_t, Q} \mathbf{x}_L + \mathbf{b}_{1, M_t, Q}) + \mathbf{b}_{2, M_t, Q}),$$

⁴ similar way for receive beamforming vectors

Contd. Predicted actions

- **Transmit powers:** The output of the last two sub-networks are:

$$a_{p_U} = \tanh(\mathbf{w}_{p_U} \mathbf{x}_L + b_{p_U}), \text{ and } a_{p_A} = \tanh(\mathbf{w}_{p_A} \mathbf{x}_L + b_{p_A}), \quad (14)$$

The output from sub-network is in the range $[-1, +1]$ which is shifted and scaled to the range $[0, P_u]$ and $[0, P_a]$ before using them in the environment

$$p_a = (a_{p_A} + 1)/2 \times P_a, \text{ and } p_u = (a_{p_U} + 1)/2 \times P_u, \quad (15)$$

where P_u and P_a are the maximum allowable transmit powers of the ULue and BS, respectively.

- Addition of Gaussian noise to the action explores the action space well which gives faster convergence

Alternative solution PerfCSI-DRL using MRC-ZF - needs CSI

- Maximize sum of UL and DL SINR by maximizing the SNR towards reception

$$\begin{aligned} \mathcal{P}_2 : \quad & \max_{\mathbf{w}_T} r_{DL} + r_{BS} \\ \text{s.t.} \quad & \|\mathbf{w}_R^{MRC} \mathbf{H}_{AA} \mathbf{w}_T\|^2 = 0, \quad \|\mathbf{w}_T\|^2 = 1, \end{aligned} \tag{16}$$

where

$$\mathbf{w}_r^{MRC} = \frac{(\mathbf{h}_{AU} + \mathbf{F}_{AI} \mathbf{\Theta}_U \mathbf{f}_{IU} + \mathbf{G}_{IA}^T \mathbf{\Theta}_D \mathbf{g}_{IU})^H}{\|\mathbf{h}_{AU} + \mathbf{F}_{AI} \mathbf{\Theta}_U \mathbf{f}_{IU} + \mathbf{G}_{IA}^T \mathbf{\Theta}_D \mathbf{g}_{IU}\|}. \tag{17}$$

- Note that the knowledge of every interferer is not available so \mathbf{g}_{IU} is not available

- We want to minimize the self-interference using Zero-Forcing. The precoder \mathbf{w}_T is in the orthogonal complement space of $\mathbf{w}_R^{MRC} \mathbf{H}_{AA}$. The orthogonal projection onto the orthogonal complement of the column space of $\mathbf{w}_R^{MRC} \mathbf{H}_{AA}$ is given by⁵

$$\Pi_{\mathbf{H}_{AA}^\dagger \mathbf{w}_R^{MRC}}^\perp = \mathbf{I}_{M_t} - \mathbf{H}_{AA}^\dagger \mathbf{w}_R^{MRC} (\mathbf{w}_R^{MRC} \mathbf{H}_{AA} \mathbf{H}_{AA}^\dagger \mathbf{w}_R^{MRC})^{-1} \mathbf{w}_R^{MRC} \mathbf{H}_{AA}.$$

- The optimal solution for transmit beamforming is

$$\mathbf{w}_T^{ZF} = \frac{\Pi_{\mathbf{H}_{AA}^\dagger \mathbf{w}_R^{MRC}}^\perp (\mathbf{h}_{DA} + \mathbf{g}_{DI} \Theta_D \mathbf{G}_{IA} + \mathbf{f}_{DI} \Theta_U \mathbf{F}_{AI}^T)^\dagger}{\|\Pi_{\mathbf{H}_{AA}^\dagger \mathbf{w}_R^{MRC}}^\perp (\mathbf{h}_{DA} + \mathbf{g}_{DI} \Theta_D \mathbf{G}_{IA} + \mathbf{f}_{DI} \Theta_U \mathbf{F}_{AI}^T)^\dagger\|}. \quad (18)$$

- This work can be extended to multiple users, too, and precoding will help us to beamform towards every user using the same antenna array.

⁵ \dagger represents conjugate transpose

Simulation setup and results

- BS situated at $(0, 0)$
- No direct path from the BS to the UL and DL users
- To promote communication, two RISs are placed at $(50, 22)$ and $(50, -22)$
- Static ULue and DLue at $(50, 20)$ and $(50, -20)$ respectively
- Maximum transmit power allowable at ULue and BS are $p_U^{max} = 50$ mW and $p_A^{max} = 1$ W
- For moving UEs, UEs move in a square area of 100 m^2 with an average speed of 1 m per time step
- BS-RIS and the RIS-user channels have LOS, so modeled as Rician channel

- For example, the channel between RIS2 and DLue is given by

$$\mathbf{g}_{DI} = \sqrt{\frac{\beta_{UI}}{1 + \beta_{UI}}} \mathbf{g}_{DI}^{LOS} + \sqrt{\frac{1}{1 + \beta_{UI}}} \mathbf{g}_{DI}^{NLOS}. \quad (19)$$

where β_{UI} is the Rician K -factor for the channels between RIS and users ⁶

- The other channels don't have an LOS, so modeled as Rayleigh
- The path loss between two points with distance d is modeled as,

$$PL(f_c, d)_{dB} = -20 \log_{10}(4\pi f_c/c) - 10\alpha \log(d/D_0), \quad (20)$$

where f_c is the carrier frequency, $D_0 = 1$ m, α is the path loss exponent.

⁶For more detail on how the channels are simulated in our setup, please refer to the paper

- The bandwidth where the system operates is 100 MHz, the noise power density is -174 dBm/Hz and the carrier frequency f_c is 3.5 GHz.
- For the DRL agent, the discounting factor $\gamma = 0.6$, buffer size $\tau = 10000$, and the learning rates of actor and critic networks are 0.0001 and 0.001 respectively
- The experiment is run for initial 50 episodes and each episode has 1000 time steps. The results are averaged over 4 independent runs
- The experiments are performed on an NVIDIA GeForce RTX 2080 Ti GPU
- The benchmark metrics used for studying the performance are UL and DL data rates with the unit bits/sec/Hz.

Our method and baseline competitor methods

- **RandPSBF**: Agent does not receive any SINR feedback from the environment; predict the RIS phases, the beamformers, and transmit powers randomly
- **OUPSBF**: Makes use of the same DRL framework as ours, except for the action-noise where it adds the Ornstein Uhlenbeck noise to the RIS phases and beamformer
- Proposed Minimum Signalling Feedback (MSF) DRL method with LSSIC and HSIC (**MSF-DRL-LSSIC** and **MSF-DRL-HSIC**):
 - The critic network and the feature extractor actor-network are feed-forward networks with two layers, each with 100 neurons.
 - At the beginning of every episode⁷, the MSF-DRL agent chooses actions randomly.
 - MSF-DRL uses a Gaussian action noise with zero mean and linearly decaying standard deviation (SD) with an initial SD of 0.3 decaying over 50 episodes.

⁷ Episode is an independent game or sequence of states where the agent and environment interacts. It starts at an initial state and ends at a terminal state.

- **PerfCSI-DRL and NoisCSI-DRL:**
 - A hypothetical experiment called “PerfCSI-DRL” to predict only the RIS phases where the perfect CSI knowledge including residual SI is available to the agent therefore serves as benchmark
 - The agent calculates the beamformers based on the ZF and MMSE principle that needs perfect CSI
 - Periodically receive the CSI to calculate the beamformers and therefore incurs overhead
 - CSI estimation methods may not be exact - NoisCSI-DRL
- **MSF-DRL-LSSIC-pos:** Along with previous actions, UL SINR and DL SINR, the past positions for a window is also given to the MSF-DRL-LSSIC agent

Comparison with static UE

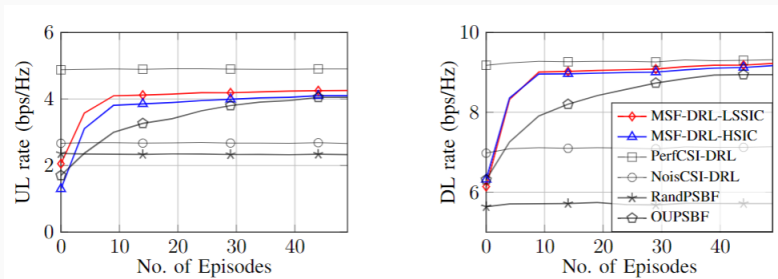


Figure 3: Rate evolution during learning in the static UE scenario. MSF-DRL-LSSIC and HSIC learns to predict and tries to reach the benchmark PerfCSI-DRL, performs much better than the case of NoisCSI-DRL.

Comparison with moving UE

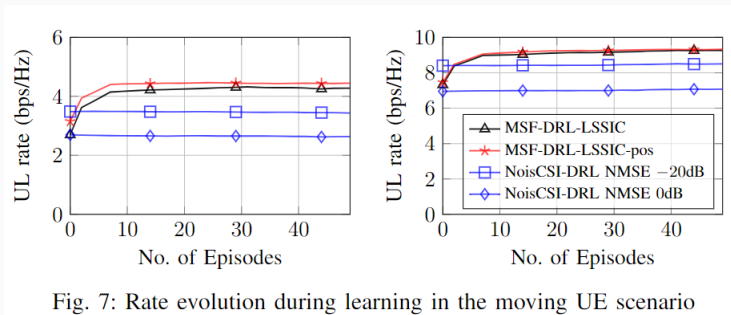


Fig. 7: Rate evolution during learning in the moving UE scenario

Figure 4: Rate evolution during learning in the moving UE scenario. In moving UE case also, the proposed method does not need any CSI and the knowledge of residual SI still performs better than NoisCSI-DRL methods.

Quantized MSF-DRL-LSSIC

- Every element of the real-valued phases is to be represented with just n bits so that instead of a real (i.e., 64 bit) phase-shift value, only n ($n \ll 64$) bit information is transmitted from the BS to the RIS
- Number of phase values that each of the passive elements can take is $Q = 2^n$ and are given by $\mathbf{p} = 2\pi/2^n \times [0, \dots, 2^n - 1]$ radians
- The architecture of the first two sub-networks for predicting RIS phases are now modified so that instead of a single phase value, they predict the probabilities of picking that phase value out of the possible values of discrete RIS phase angles of length $Q = 2^n$
- The sub-network for predicting the phase-shift of RIS1 takes the feature \mathbf{x}_L as input and passes it through N_1 fully connected layers, each with 2^n neurons, followed by a softmax activation on each of the N_1 outputs.

Results with quantized phase MSF-DRL-LSSIC

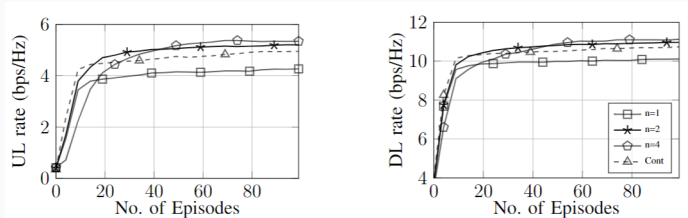


Figure 5: Effect of quantization on phaseshifts of RISs on MSF(Q)-DRL methods. When $n = 1$, the phase-shifts can take $Q = 2^n = 2$ values $\{0, \pi\}$; for $n = 0$, the phase-shifts are fixed to 0 and only the beamformers are learned.

Reduced signalling with two-stage DRL algorithm

- Uses only one **scalar feedback**, i.e., the weighted sum of UL and DL SINR from the environment instead of CSI which is in the form of Matrix for MIMO channels
- For our simulation setup i.e. 10 BS antenna elements and 36 RIS elements per RIS, the number of instantaneous CSI that need to be fed to the algorithm is 813
- Proposed MSF-DRL method needs to transmit one pilot signal and needs to receive one reward values in a single time step

Table 1: Proposed methods reduces signaling

| Methods | FLOPs | Needs CSI, SI (#values) | Signaling from BS to RIS (in bits) |
|-------------|--------------------|-------------------------|------------------------------------|
| PerfCSI-DRL | 2.9×10^4 | Yes (813) | $64N_1 + 64N_2 = 4608$ |
| MSF-DRL | 3.3×10^4 | No (2) | $64N_1 + 64N_2 = 4608$ |
| MSF(Q)-DRL | 5.46×10^4 | No (2) | $2N_1 + 2N_2 = \mathbf{144}$ |

- The two-stage learning algorithm **assumes the absence of a good SI mitigation scheme** and the costly CSI overhead but **still performs almost as well as perfect CSI-based semi-oracle DRL methods**.
- To overcome the challenge of SI, a **least square-based method** is proposed when a good estimate of SI is not present.
- The performances are shown in scenario with moving UEs as well
- We also propose a DRL framework that **can learn quantized RIS phase shifts**. The quantized phase DRL method has 32 times lesser signaling than the continuous phase DRL method with better convergence.

Thank you!

- To get stable, uncorrelated gradients for policy improvement, DDPG maintains a replay buffer of finite size τ and samples the observations from the buffer in mini-batches to update the parameters.
- At each timestep, the state $\mathbf{s}^{\{i\}}$ and the action taken $\mathbf{a}^{\{i\}}$ along with the reward obtained $r^{\{i\}}$ and the next state $\mathbf{s}^{\{i+1\}}$ is stored as an experience $(\mathbf{s}^{\{i\}}, \mathbf{a}^{\{i\}}, r^{\{i\}}, \mathbf{s}^{\{i+1\}})$ to the buffer \mathcal{B} .
- DDPG also uses target networks with parameters $\bar{\omega}^a$ and $\bar{\omega}^c$ to avoid divergence in value estimation
- For the critic network $\mathbb{C}(\cdot|\omega^c)$ to compute the Q-value for each state action-pair, an estimate of return for state s_i in each sample is computed as

$$y^{\{i\}} = r^{\{i\}} + \gamma \mathbb{C}(\mathbf{s}^{\{i+1\}}, \mathbb{A}(\mathbf{s}^{\{i+1\}}|\bar{\omega}^a)|\bar{\omega}^c). \quad (21)$$

- Once we observe a reward r_i after taking an action a_i , based on the estimate for return, the mean squared Bellman error (MSBE) is computed as

$$\mathcal{L} = \frac{1}{N} \sum_i \left(y^{\{i\}} - \mathbb{C}(\mathbf{s}^{\{i\}}, \mathbf{a}^{\{i\}} | \boldsymbol{\omega}^c) \right)^2, \quad (22)$$

where, $\mathbb{C}(\cdot)$ is the predicted output value of critic network with parameter $\boldsymbol{\omega}^c$ for the state $\mathbf{s}^{\{i\}}$ and action $\mathbf{a}^{\{i\}}$ before seeing the reward.

- Then, the critic network parameters are updated as

$$\boldsymbol{\omega}^c \leftarrow \boldsymbol{\omega}^c - \eta_c \nabla_{\boldsymbol{\omega}^c} \mathcal{L}, \quad (23)$$

where $\eta_c \ll 1$ is the stepsize for the stochastic update.

- In our case the critic network is a **fully connected network** which gives a scalar output as Q-value

- For the actor-network, the update depends on both the gradient of action as well as the improvement in Q-value. The final update for updating parameters of actor-network ω^a is given by

$$\omega^a \leftarrow \omega^a + \eta_a \frac{1}{N} \sum_i (\nabla_{\omega^a} \mathbb{A}(s) \nabla_a \mathbb{C}(s, a)|_{a=\mathbb{A}(s)}) , \quad (24)$$

where $\eta_a \ll 1$ is the update stepsize.

- Finally, the target network parameters are updated in every U timestep to provide stable value estimates using an exponentially weighted update as $\bar{\omega}^c \leftarrow \lambda \omega^c + (1 - \lambda) \bar{\omega}^c$, and $\bar{\omega}^a \leftarrow \lambda \omega^a + (1 - \lambda) \bar{\omega}^a$, with $\lambda \ll 1$.