

# Multi-agent deep reinforcement learning (MADRL) meets multi-user MIMO systems

Heunchul Lee  
Ericsson Research,  
Ericsson AB  
Stockholm, Sweden  
heunchul.lee@ericsson.com

Jaeseong Jeong  
Ericsson Research,  
Ericsson AB  
Stockholm, Sweden  
jaeseong.jeong@ericsson.com

**Abstract**—A multi-agent deep reinforcement learning (MADRL) is a promising approach to challenging problems in wireless environments involving multiple decision-makers (or actors) with high-dimensional continuous action space. In this paper, we present a MADRL-based approach that can jointly optimize precoders to achieve the outer-boundary, called pareto-boundary, of the achievable rate region for a multiple-input single-output (MISO) interference channel (IFC). In order to address two main challenges, namely, multiple actors (or agents) with partial observability and multi-dimensional continuous action space in MISO IFC setup, we adopt a multi-agent deep deterministic policy gradient (MA-DDPG) framework in which decentralized actors with partial observability can learn a multi-dimensional continuous policy in a centralized manner with the aid of shared critic with global information. Meanwhile, we will also address a phase ambiguity issue with the conventional complex baseband representation of signals widely used in radio communications. In order to mitigate the impact of phase ambiguity on training performance, we propose a training method, called phase ambiguity elimination (PAE), that leads to faster learning and better performance of MA-DDPG in wireless communication systems. The simulation results exhibit that MA-DDPG is capable of learning a near-optimal precoding strategy in a MISO IFC environment.

**Index Terms**—Multi-agent deep reinforcement learning (MADRL), Multi-agent deep deterministic policy gradient (MA-DDPG), Multiple-input multiple-output (MIMO), Interference Channel (IFC)

## I. INTRODUCTION

### A. Multi-cell MIMO problems and multi-agent system

As cellular data demand continues to rise, an ultra-dense network is widely considered as a key component in managing this rising. Multiple-input multiple-output (MIMO) technique has been developed for efficient transmission and reception of radio signals in multiple antenna systems. In particular, downlink multi-user MIMO is a promising technique to achieve higher throughput in a multi-cell environment. However, in general, the optimization problems in a multi-cell multi-user MIMO system are nonconvex and difficult to solve using the traditional approach based on mathematical models. Machine learning (ML) is a promising approach to overcome the limitations of the traditional model-based approach, allowing the future cellular networks to evolve towards more scalable and intelligent architectures [1]. In this paper, by leveraging

the recent success of deep reinforcement learning (DRL) and multi-agent (MA) learning [2] [3] [4], we propose a ML-based approach that integrates multi-agent deep reinforcement learning (MADRL) into downlink multi-cell multi-user wireless systems.

Ideally, joint data transmission schemes assume full MIMO cooperation in a multi-cell multi-user environment. For instance, coordinated multi-point (CoMP) with joint transmission (JT) is a cellular data transmission technique involving simultaneous transmission from multiple base stations (BSs) to the same user [5]. However, potential solutions to the JT scheme require significant amounts of global channel state information (CSI) and data sharing between the base stations, which is not only expensive but also difficult in real-world cellular systems. The optimal JT scheme can be reduced to coordinated beamforming (CB) schemes based on the transmission of the signal by a single base station that require local CSI and no inter-cell data sharing [5]. In this case, each BS operates independently by treating the interference as background noise and the multi-cell multi-user setup can be modeled as MIMO interference channel (IFC). Compared to a single-cell system, the performance of MIMO IFC can be severely impacted by inter-cell interference, which becomes a crucial limiting factor.

Reinforcement learning (RL) allows an agent to learn the optimal action policy that returns the maximum reward through trial-and-error interactions with a challenging dynamic environment [6]. RL has been used to solve challenging problems in various areas ranging from games to robotics. In wireless, RL is also emerging as one of key enablers for designing 6G AI-driven PHY-layer [1]. Recently, we have investigated RL-based approaches to improve the performance of MIMO systems [7] [8]. However, these studies have focused on enhancing the performance of single-cell MIMO systems. Multi-cell multi-user precoding problems can be seen as a multi-agent system that learns to coordinate transmission schemes (or action policies) in interaction with other base stations (or other agents). Therefore, scaling our previous work to more complex multi-agent problems is crucial to building future intelligent networks that can operate in real-world multi-cell environments.

The multi-agent problem requires complex inter-cell interference coordination in the sense that each BS should exhibit cooperative behavior to maximize the signal power to a desired user while minimizing the interference power to other users in the multi-cell environment. We note that this problem setup poses two main challenges: i) *multiple actors (or agents) with partial observability* and ii) *multi-dimensional continuous action space*. The first challenge is a direct result of practical limitations of accessible information by local agents distributed in MISO-IFC, and the second challenge comes from the fact that multi-dimensional precoding vectors should be optimized for multi-antenna BSs based on a certain transmit power constraint.

### B. Main contributions

To address these two challenges, we propose a multi-agent deep deterministic policy gradient (MA-DDPG)-based approach that can learn an optimal precoding strategy in multi-cell multi-user MIMO systems under the assumptions of local CSI and no inter-cell data sharing. In particular, in order to permit tractable performance analysis, we consider a multiple-input single-output (MISO) IFC in which two base stations equipped with multiple antennas serve two single-antenna users, making the precoding problem tractable by the numerical methods proposed in [9] and [10]. In this two-user MISO IFC setup, we can obtain the achievable rate region by using the work of [9] and derive the pareto-boundary of the rate region.

MA-DDPG algorithms provide a multi-agent framework to learn a high dimensional continuous policy [3] [4]. The actor-critic based policy gradient algorithms allow centralized training with decentralized execution in which local actors with partial observability can learn a globally optimal policy with the aid of centralized critic with global information at training time and execute the learned policy based only on partial observations at execution time. At the same time, the deterministic policy gradient (DPG) algorithm enables the agents to learn a multi-dimensional continuous policy. The MA-DDPG framework is adopted to improve the quality of the received signals in MISO IFC by alleviating the inter-cell interference. Meanwhile, we also investigate the impact of phase ambiguity with the baseband representation of wireless channel on training performance. The complex-valued representation of channel states has inherent phase ambiguity in the sense that the phase-shifted versions of a channel state will have the same impact in system performance as the original channel state. From the wireless system design point of view, this phase ambiguity should not be a problem but it can cause a performance degradation in a multi-agent learning system. In order to mitigate the impact of this phase ambiguity in training time, we propose a feature engineering method, called *phase ambiguity elimination (PAE)*, as a pre-processing step on input channel states to using a MA-DDPG algorithm to learn an optimal policy. By applying the proposed PAE method in the state space, we demonstrate faster learning and better performance

of MA-DDPG in MISO IFC. The simulation results indicate that MA-DDPG is capable of learning precoding schemes which achieves the outer boundary, called pareto-boundary, of achievable rate regions on MISO IFC environments. To the best of our knowledge, this is the first work to demonstrate that the MA-DDPG framework can jointly optimize precoders to achieve the pareto-boundary of achievable rate region in a multi-cell multi-user multi-antenna system.

## II. SYSTEM MODEL

In this section, we describe the system model of MISO IFC, where two BSs equipped with multiple antennas simultaneously communicate with its own desired user equipment (UE) equipped with a single antenna in the same time-frequency resource. It is important to recall that we have assumed this MISO IFC setup to ensure that the rate region can be obtained by the numerical method in [9]. We also describe the numerical method for obtaining achievable rate region as well as two pareto-optimal rate pairs with closed-form expressions. In the simulation section, numerical results will be used as a quantitative criterion for demonstrating the optimality of MA-DDPG in multi-cell multi-user MISO systems.

### A. MISO IFC scenario

In this subsection, we present a MISO IFC model and related assumptions. As shown in Figure 1, BS  $i \in \{1, 2\}$  desires to send the data symbol  $d_i$  to UE  $i$ . The base stations employ  $n_t$  transmit antennas and each UE is equipped with a single receive antenna. BS  $i$  employs a linear precoding vector  $\mathbf{w}_i$  of size  $n_t$ -by-1 prior to transmission over the air, which transforms the data symbol  $d_i$  to the  $n_t$ -by-1 transmitted vector  $\mathbf{x}_i = \mathbf{w}_i d_i$ . The channel model from the BS  $i$  to the two UEs are represented by an 1-by- $n_t$  channel vector  $\mathbf{h}_i = [h_{i,1}, h_{i,2}, \dots, h_{i,n_t}]$  and  $\mathbf{g}_i = [g_{i,1}, g_{i,2}, \dots, g_{i,n_t}]$ , where the  $j$ -th elements  $h_{i,j}$  and  $g_{i,j}$  denote the path gain from the  $j$ -th antenna of BS  $i$  to the desired UE  $i$  and the other UE, respectively. The channel elements are independently and identically distributed (i.i.d.) according to  $\mathcal{N}_{\mathbb{C}}(0, 1)$ , i.e.,  $\mathbf{h}_i \in \mathbb{C}^{n_t}$  and  $\mathbf{g}_i \in \mathbb{C}^{n_t}$ .

The received signal  $y_i$  at UE  $i$  can be expressed as, for  $i = 1$ ,

$$y_1 = \mathbf{h}_1 \mathbf{x}_1 + \mathbf{g}_2 \mathbf{x}_2 + n_1 = \mathbf{h}_1 \mathbf{w}_1 d_1 + \mathbf{g}_2 \mathbf{w}_2 d_2 + n_1, \quad (1)$$

and, for  $i = 2$ ,

$$y_2 = \mathbf{h}_2 \mathbf{x}_2 + \mathbf{g}_1 \mathbf{x}_1 + n_2 = \mathbf{h}_2 \mathbf{w}_2 d_2 + \mathbf{g}_1 \mathbf{w}_1 d_1 + n_2, \quad (2)$$

where  $n_i$  denotes complex-valued additive white Gaussian noise (AWGN) at UE  $i$ , distributed as  $\mathcal{N}_{\mathbb{C}}(0, \sigma_n^2)$ . Note that UE  $i$  not only receives its desired signal  $\mathbf{w}_i d_i$  through the channel  $\mathbf{h}_i$  but also the inter-cell interference (ICI)  $\mathbf{w}_j s_j$  from the signal intended for the other UE  $j \neq i$ . We impose the power constraint  $\mathbb{E}[\text{tr}[\mathbf{x}_i \mathbf{x}_i^H]] = 1$  under assumption of unit-norm weight vectors  $\mathbf{w}_i$  and unit-power symbols  $d_i$ , i.e.,  $\|\mathbf{w}_i\| = 1$  and  $\mathbb{E}[|d_i|] = \sigma_d^2 = 1$ , where  $\mathbb{E}[\cdot]$  denotes the expectation with respect to the distribution of the underlying random variable,

$\text{tr}[\cdot]$  denotes the trace operator of a matrix,  $\|\cdot\|$  indicates the 2-norm of a vector, and  $|\cdot|$  denotes the absolute value of a scalar. Then the average transmit signal-to-noise ratio (SNR) of the network is defined as  $\rho = \frac{1}{\sigma_n^2}$ .

In the multi-cell environment, the signal quality is measured in the form of achievable data rate as a function of received signal-to-interference-plus-noise ratio (SINR). From (1), the received SINR for UE 1 is defined as

$$\text{SINR}_1 = \frac{\sigma_d^2 |\mathbf{h}_1 \mathbf{w}_1|^2}{\sigma_n^2 + \sigma_d^2 |\mathbf{g}_2 \mathbf{w}_2|^2}, \quad (3)$$

Similarly, the received SINR at UE 2 is given by

$$\text{SINR}_2 = \frac{\sigma_d^2 |\mathbf{h}_2 \mathbf{w}_2|^2}{\sigma_n^2 + \sigma_d^2 |\mathbf{g}_1 \mathbf{w}_1|^2}. \quad (4)$$

### B. Rate region and two pareto-optimal rate pairs with closed-form expressions

Let  $r_i$  denote the rate for UE  $i$ . We denote by  $R(\mathbf{w}_1, \mathbf{w}_2)$  the conditional rate tuple  $(r_1, r_2)$  that can be achieved for a given pair of linear precoding vectors  $\mathbf{w}_1$  and  $\mathbf{w}_2$ . The theoretical limit of rate tuple  $R(\mathbf{w}_1, \mathbf{w}_2)$  achievable with Gaussian random coding is given as

$$R(\mathbf{w}_1, \mathbf{w}_2) = \left[ \log_2 \left( 1 + \frac{\sigma_d^2 |\mathbf{h}_1 \mathbf{w}_1|^2}{\sigma_n^2 + \sigma_d^2 |\mathbf{g}_2 \mathbf{w}_2|^2} \right), \log_2 \left( 1 + \frac{\sigma_d^2 |\mathbf{h}_2 \mathbf{w}_2|^2}{\sigma_n^2 + \sigma_d^2 |\mathbf{g}_1 \mathbf{w}_1|^2} \right) \right]. \quad (5)$$

Then, the achievable rate region  $\mathcal{R}$  can be defined as the closure of the set of all achievable rate pairs  $R(\mathbf{w}_1, \mathbf{w}_2)$  under the power constraints  $\|\mathbf{w}_i\|^2 \leq 1$  for  $i = 1, 2$ ,

$$\mathcal{R} = \cup_{\|\mathbf{w}_i\|^2 \leq 1, i=1,2} R(\mathbf{w}_1, \mathbf{w}_2). \quad (6)$$

Our goal is to find a pareto-optimal linear precoding scheme to construct pairs of precoding vectors  $\mathbf{w}_1$  and  $\mathbf{w}_2$  that achieve all the rate pairs on the pareto-boundary of the achievable rate region  $\mathcal{R}$ . In general, the rate region is not known. But, we can numerically obtain the pareto-boundary by using the work in [9]. As shown in [9], any pareto-optimal rate pairs can be achieved by using precoding vectors  $\mathbf{w}_1$  and  $\mathbf{w}_2$  that are parameterized by maximum ratio transmission (MRT) and zero-forcing (ZF) solution. The MRT and ZF solution at BS  $i$  are given by

$$\mathbf{w}_i^{\text{mrt}} = \mathbf{h}_i^H. \quad (7)$$

and

$$\mathbf{w}_i^{\text{zf}} = (\mathbf{g}_i^H \mathbf{g}_i)^{-1} \mathbf{h}_i^H. \quad (8)$$

The MRT solution  $\mathbf{w}_i^{\text{mrt}}$  is optimal for single-user MIMO by maximizing the signal gain at the intended UE  $i$ . In comparison, the ZF solution  $\mathbf{w}_i^{\text{zf}}$  can be seen as a MRT precoding vector designed for projecting the symbol  $d_i$  on the null space of  $\mathbf{g}_i$ .

As shown in Figure 3, we can numerically evaluate the rate region  $\mathcal{R}$  in (6) by using the parameterized precoding vectors  $\mathbf{w}_1$  and  $\mathbf{w}_2$  presented in [9], and determine the pareto-boundary for two-user MISO IFC.

We now briefly present two pareto-optimal rate pairs with closed-form expressions. The first reference rate pair can be directly achieved by applying the mixed pairs of two precoding vectors given in (7) and (8), resulting in two rate pairs  $R(\mathbf{w}_1^{\text{mrt}}, \mathbf{w}_2^{\text{zf}})$  and  $R(\mathbf{w}_1^{\text{zf}}, \mathbf{w}_2^{\text{mrt}})$ . As will be seen later in Figure 4, the two rate pairs correspond to two corner points on the pareto-boundary.

The second reference rate pair is given by the leakage-based precoder scheme proposed in [11] that can maximize the received SINR at the receiver by maximizing the signal-to-leakage-and-noise ratio (SLNR) at the transmitter.

The SLNR at the UE 1 is defined as

$$\text{SLNR}_1 = \frac{\sigma_d^2 |\mathbf{h}_1 \mathbf{w}_1|^2}{\sigma_n^2 + \sigma_d^2 |\mathbf{g}_1 \mathbf{w}_1|^2}, \quad (9)$$

where the second term on the denominator  $\mathbf{g}_1 \mathbf{w}_1$  indicates the *leakage* caused by the signal intended for the desired UE 1 to the other UE 2.

The optimal SLNR solution can be obtained by

$$\mathbf{w}_1^{\text{slnr}} = \text{max-eigenvector} \left( (\sigma_n^2 \mathbf{I} + \mathbf{g}_1^H \mathbf{g}_1)^{-1} \mathbf{h}_1^H \mathbf{h}_1 \right). \quad (10)$$

Similarly, the optimal SLNR solution at UE 2 is given by

$$\mathbf{w}_2^{\text{slnr}} = \text{max-eigenvector} \left( (\sigma_n^2 \mathbf{I} + \mathbf{g}_2^H \mathbf{g}_2)^{-1} \mathbf{h}_2^H \mathbf{h}_2 \right). \quad (11)$$

The SLNR solution is known to achieve a sum-rate point on the pareto-boundary, corresponding to the rate pair in rate region that obtains the maximum sum rate. This sum-rate optimal rate pair, denoted by  $R(\mathbf{w}_1^{\text{slnr}}, \mathbf{w}_2^{\text{slnr}})$ , as well as the two corner points  $R(\mathbf{w}_1^{\text{mrt}}, \mathbf{w}_2^{\text{zf}})$  and  $R(\mathbf{w}_1^{\text{zf}}, \mathbf{w}_2^{\text{mrt}})$  will serve as an upper reference rate pair when we provide simulation results of MA-DDPG in the simulation section.

### III. PARETO-OPTIMAL PRECODING STRATEGY

In this section, we present a pareto-optimal precoding strategy based on MA-DDPG framework for a MISO IFC setup. We first describe how MA-DDPG framework can be adopted to learn an optimal precoding strategy in MISO IFC setup. Then, we address the phase ambiguity issue with the conventional complex baseband representation in radio communications. In order to avoid the impact of this phase ambiguity in the learning process, we propose a training method that leads to faster learning and better performance of MA-DDPG in wireless communication systems.

#### A. Multi-agent RL with a continuous action space

RL problems can be formalized by modelling the interaction between the agent and the environment as a *Markov decision process* (MDP). An MDP consists of a set of environment states  $\mathcal{S}$ , a set of available actions  $\mathcal{A}$ , a reward  $r \in \mathbb{R}$  and a state transition function  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  from one state to

another given an action taken. A policy is a mapping function from state to action in the MDP that specifies action  $a$  that is taken in state  $s$ . At each time step, an agent observes a state  $s \in \mathcal{S}$  and chooses an action  $a \in \mathcal{A}$ . After each time step, the agent gets an immediate reward  $r$  and next state  $s' \in \mathcal{S}$  in return for the action taken. In this paper, the policy is assumed to be a deterministic function, denoted by  $a = \mu(s)$ .

DPG algorithm can be used to handle multi-dimensional continuous actions in many continuous control problems [2]. DPG utilizes a novel actor-critic architecture that consists of two components, namely, actor and critic [6]. The actor learns to produce a deterministic policy based on the state, while the critic learns to estimate the true action-value (Q-value) function of an action given a state. The policy and the value function are parameterized by neural networks as  $\mu_\phi$  and  $Q_\theta^\mu$ , respectively. The critic estimates the policy gradient from the learned  $Q_\theta^\mu$  and sends it to the actor to update the policy  $\mu_\phi$  at the same time. As a deep variant of DPG, DDPG combines deep neural networks with the actor-critic architecture [2].

MA-DDPG is increasingly used within a diverse range of applications involving multi-agent environments with multi-dimensional continuous action space. There are basically two design factors to consider: allowing agents to share a single critic and augmenting a critic with the policies of other agents. By combining features of the original designs in [3] and [4], in this paper we consider a MA-DDPG framework where one critic is shared by all agents and the centralized critic is augmented with policies of the agents.

### B. MA-DDPG in MISO IFC

Figure 1 illustrates MA-DDPG model for the MISO IFC setup under the following assumptions:

- Non-real time communication available for the critic to learn an action-value function in a centralized manner based on global information about environmental channel states and policies of both agents.
- To fulfill the real or near-real time requirement of precoding scheme, each agent  $i$  at the  $i$ -th BS chooses a precoding vector  $\mathbf{w}_i$  based on local information given by the partial channel observation  $\mathbf{h}_i$  and  $\mathbf{g}_i$  only.

As shown in Figure 1, the MA-DDPG extends the actor-critic policy gradient method to provide a framework of centralized training with decentralized execution, improving training stability at training phase and performance robustness at execution phase. More importantly, the framework allows non-real time learning based on global information and real or near-real time execution based on local information, therefore making it practical for real-world cellular environments.

In TDD-based systems, downlink CSI can be derived from uplink channel observations thanks to channel reciprocity. For instance, in current LTE and NR specifications, each BS  $i$  can estimate downlink channel state of  $\mathbf{h}_i$  and  $\mathbf{g}_i$  based on uplink sounding reference signal (SRS) transmitted by the desired and interfering UE. Therefore, the partial state  $\mathbf{s}_i$  (or

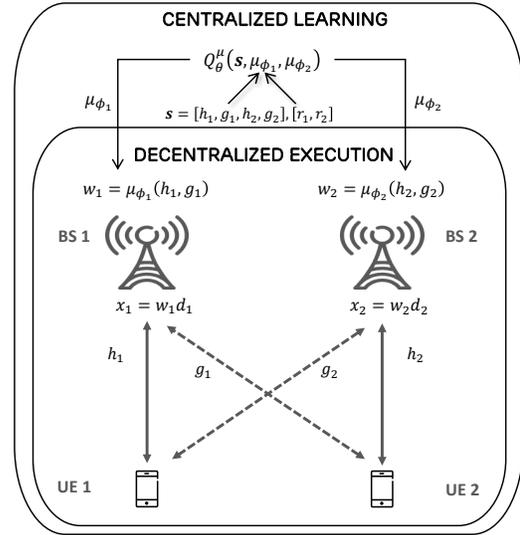


Fig. 1. MA-DDPG in MISO IFC

partial observation  $\mathbf{o}_i$  more specifically) at BS  $i$  can be defined, directly from the local observation of  $\mathbf{h}_i$  and  $\mathbf{g}_i$ , as

$$\mathbf{s}_i = [\mathbf{h}_i, \mathbf{g}_i]. \quad (12)$$

The actor  $i$  chooses an action for a given state  $\mathbf{s}_i$  by using a deterministic policy

$$\mathbf{a}_i = \mu_{\phi_i}(\mathbf{s}_i), \quad (13)$$

where the output action  $\mathbf{a}_i = [a_{i,1}, a_{i,2}, \dots, a_{i,n_t}]^T$  is used to determine a precoding vector  $\mathbf{w}_i = [w_{i,1}, w_{i,2}, \dots, w_{i,n_t}]^T$ .

While we use a deterministic policy  $\mu_{\phi_i}$  that always yields the same action for the same state, a stochastic policy is desirable for exploration at training time. DDPG uses a stochastic behavior policy to select actions, different from the learned policy. In order to ensure sufficient exploration, in this paper, we perturb the deterministic action  $a = \mu_{\phi_i}(s)$  by adding a noise vector whose entries are i.i.d. according to  $\mathcal{N}_{\mathbb{C}}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ , as described in [2].

Each agent  $i$  receives as a reward the rate  $r_i$  given by (5) as a function of the environmental state in (12) and actions taken by the stochastic behavior policy. In order to achieve a rate pair on the pareto-boundary of the achievable rate region, the agents should behave in a cooperative manner to maximize the collective reward. The collective reward  $r_{c,\alpha}$  for achieving a pareto-optimal precoding strategy can be defined as a weighted sum of achieved rates  $r_1$  and  $r_2$  [10]

$$r_{c,\alpha} = \alpha \cdot r_1 + (1 - \alpha)r_2 \quad (14)$$

where  $\alpha \in [0, 1]$  denotes the weighting scalar between the two rewards  $r_1$  and  $r_2$ .

We note that the weighting factor  $\alpha$  determines which pareto-optimal rate pair to achieve by specifying a straight line with slope  $-\frac{\alpha}{(1-\alpha)}$  to the pareto-boundary. In the simulation section, we will consider three different values,  $\alpha = 1/2, 2/3,$

and  $3/4$ , which will achieve the point of rate pairs on the pareto-boundary where the slope of tangent is equal to  $-1$ ,  $-2$  and  $-3$ , respectively.

The centralized critic aims to maximize its total expected return  $R_c = \sum_{t=0}^{\infty} \gamma^t r_{c,\alpha}^t$ , where  $\gamma \in [0, 1]$  denotes a discounting factor to the future rewards  $r_{c,\alpha}^t$  for  $t = 1, 2, \dots$ , relative to the immediate reward  $r_{c,\alpha}^0$ . To this end, the true action-value function is approximated by the critic network

$$\mathbf{Q}_{\theta}^{\mu}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) \quad (15)$$

where  $\mu = [\mu_{\phi_1}, \mu_{\phi_2}]$  and the global state  $\mathbf{s}$  is given by  $\mathbf{s} = [\mathbf{s}_1, \mathbf{s}_2]$ .

Equation (15) shows that the critic function uses the policies of both agents so that each agent can learn approximate models of other agents from the learned critic.

At each time step, according to Q-learning, the critic updates the value-function parameters  $\theta$  as follows:

$$\theta \leftarrow \theta + \eta_c (Y^{\theta} - \mathbf{Q}_{\theta}^{\mu}(\mathbf{s}, \mathbf{a}_1^p, \mathbf{a}_2^p)) \nabla_{\theta} \mathbf{Q}_{\theta}^{\mu}(\mathbf{s}, \mathbf{a}_1^p, \mathbf{a}_2^p), \quad (16)$$

where  $\eta_c \in [0, 1]$  is a critic learning rate,  $\nabla_{\theta} \mathbf{Q}_{\theta}^{\mu}(\cdot)$  denotes the vector of partial derivatives with respect to the components of  $\theta$ ,  $\mathbf{a}_i^p$  is the noisy version of  $\mathbf{a}_i = \mu_{\phi_i}(\mathbf{s}_i)$ , and  $Y^{\theta}$  indicates the newly estimated value on the current step, assuming the deterministic actions on the next step  $\mathbf{s}' = [\mathbf{s}'_1, \mathbf{s}'_2]$ , which is given by

$$Y^{\theta} = r_{c,\alpha} + \gamma \mathbf{Q}_{\theta}^{\mu}(\mathbf{s}', \mu_{\phi_1}(\mathbf{s}'_1), \mu_{\phi_2}(\mathbf{s}'_2)). \quad (17)$$

The actor  $i$  estimates the policy parameters  $\phi_i$  that maximize the expected reward by updating the parameters via a gradient ascent

$$\phi_i \leftarrow \phi_i + \eta_a \nabla_{\phi_i} J(\mu_{\phi_i}), \quad (18)$$

where  $\eta_a$  is an actor learning rate and, according to the DPG algorithm, the gradient  $\nabla_{\phi_i} J(\mu_{\phi_i})$  is obtained by

$$\nabla_{\phi_i} J(\mu_{\phi_i}) = \nabla_{\phi_i} \mu_{\phi_i}(\mathbf{s}_i) \nabla_{\mathbf{a}_i} \mathbf{Q}_{\theta}^{\mu}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) |_{\mathbf{a}_i = \mu_{\phi_i}(\mathbf{s}_i)}. \quad (19)$$

After training is completed, the local actors can execute the learned policies only based on the local CSI at execution phase, successfully addressing the challenge of multiple agents with partial observability.

### C. Phase ambiguity elimination

In this subsequent section, we address the phase ambiguity issue in the commonly used vector (or matrix) representation of wireless channel states and then present a method to improve MA-DDPG training in MISO IFC.

In radio communications, a passband channel state is represented by a complex-valued baseband equivalent. In our signal model given in (1) and (2), the channel vectors  $\mathbf{h}_1$ ,  $\mathbf{g}_1$ ,  $\mathbf{h}_2$  and  $\mathbf{g}_2$  can be expressed as a complex-valued representation with respect to amplitude and phase, denoted by

$$\mathbf{c} = [a_1 \exp(j\vartheta_1), a_2 \exp(j\vartheta_2), \dots, a_{n_t} \exp(j\vartheta_{n_t})], \quad (20)$$

where  $a_i$  and  $\vartheta_i$  are the amplitude and phase of the  $i$ -th element of vector  $\mathbf{c}$ .

We note that the channel state given by  $\mathbf{c}$  in (20) has inherent phase ambiguity resulting from the complex-valued baseband signal representation. More specifically, all the phase-shifted states  $\exp(j\varphi)\mathbf{c}$  with arbitrary phases  $\varphi$  are supposed to lead to the same action as that of the original state  $\mathbf{c}$ . From the wireless system design point of view, this phase ambiguity should not be a problem, but it introduces a many-to-one mapping issue between a set of phase-shifted states with different offsets and one target optimal action, which further complicates the training task, and thereby, degrades the performance in a multi-agent learning system. In order to combat the degradation due to the many-to-one mapping nature of state-action pairs, we propose a PAE method as a pre-processing on channel states  $\mathbf{h}_1$ ,  $\mathbf{g}_1$ ,  $\mathbf{h}_2$  and  $\mathbf{g}_2$ , that maps channel states with phase ambiguity into the same state. The training method can utilize any mapping function  $f_{PAE}(\cdot)$  that eliminates inherent phase ambiguity. For instance, in this paper we consider a PAE mapping function that maps each state  $\mathbf{c}$  onto one state whose first element is purely real-valued as follows:

$$f_{PAE}(\mathbf{c}) = [a_1, a_2 \exp(j(\vartheta_2 - \vartheta_1)), \dots, a_{n_t} \exp(j(\vartheta_{n_t} - \vartheta_1))]. \quad (21)$$

In summary, the final state representation can be obtained by applying the mapping function  $f_{PAE}$  to  $\mathbf{h}_i$  and  $\mathbf{g}_i$  as

$$\mathbf{s}_i^{PAE} = [f_{PAE}(\mathbf{h}_i), f_{PAE}(\mathbf{g}_i)]. \quad (22)$$

In the following section, we will show that the proposed PAE training method can achieve a faster convergence and a better performance compared to the trivial approach given in (12) without consideration of phase ambiguity elimination.

## IV. NUMERICAL RESULTS

In this section, we provide simulation results and comparisons with the pareto-boundary to demonstrate the optimality of MA-DDPG framework in MISO IFC setup. We consider BSs with three transmit antennas, i.e.,  $n_t = 3$ , and fix the average SNR to be 10 dB in all the simulations. Note that the all the precoding weight vectors given by the numerical methods and learned by MA-DDPG should be normalized to make  $\|\mathbf{w}_i\| = 1$ . We implemented the MA-DDPG model in TensorFlow 2, training a critic and two actors all with three hidden fully-connected layers. Each episode length is defined to have 10,000 time steps. We start from a perturbation variance  $\sigma_p^2 = 0.1$  and multiply it by a decaying factor of 0.993 over episodes.

Figure 2 illustrates the impact of phase ambiguity on learning performance, comparing with the theoretical upper-bound of MA-DDPG with  $\alpha = 1/2$ . After each episode, consisting of 10,000 training steps, we evaluate the learned policies in terms of average sum rate using a test set of 5000 unseen samples. MA-DDPG with PAE achieves more than 99% of the maximum achievable sum rate after 97 episodes

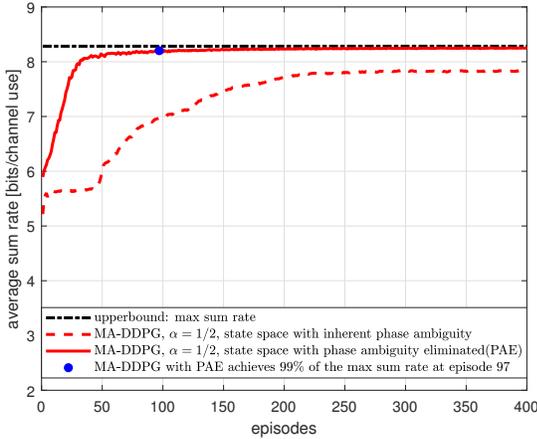


Fig. 2. Training convergence and performance impact due to phase ambiguity

while only a maximum of 94% is achieved without PAE due to the many-to-one mapping nature of state-action pairs. The simulation results show that the proposed PAE training method can achieve a faster and even better learning curve in a MISO IFC environment.

In the following figures, we provide numerical results to demonstrate the optimality of MA-DDPG with the proposed PAE method. MA-DDPG models are trained over 200 episodes and quantitative comparisons are provided in terms of achieved rate pair for a random test sample or in terms of average rate pair for 5000 test samples. For a given sample  $\mathbf{h}_1 = [-0.569 + j0.227, -0.018 + j0.456, -0.213 + j0.254]$ ,  $\mathbf{g}_1 = [-0.054 - j0.240, 0.298 - j0.232, 0.334 - j0.403]$ ,  $\mathbf{h}_2 = [-0.846 - j0.287, -0.129 + j0.073, -0.098 + j0.499]$ , and  $\mathbf{g}_2 = [0.636 - j0.493, -0.167 - j0.050, 0.204 + j0.460]$ , the achieved rate pairs by MA-DDPGs are shown in Figure 3 in comparison to the achievable rate pairs by the numerical method. Figure 4 shows learning behaviors of MA-DDPGs over 200 episodes in comparison to the pareto-optimal rate pairs, namely, the sum-rate optimal reference point  $R(\mathbf{w}_1^{slnr}, \mathbf{w}_2^{slnr})$  as well as the two corner points  $R(\mathbf{w}_1^{mrt}, \mathbf{w}_2^{zf})$  and  $R(\mathbf{w}_1^{zf}, \mathbf{w}_2^{mrt})$ , based on 5000 test samples that haven't seen by the agents before. These simulation results demonstrate that the decentralized actors with partial observability are able to discover optimal coordination strategies with the aid of the centralized critic in MISO IFC environments.

## V. CONCLUSION

In this paper, we have proposed a multi-agent deep reinforcement learning approach for precoding method in multi-cell multi-user MIMO systems. We have demonstrated that a MA-DDPG framework is able to automatically learn a near-optimal precoding policy in MISO IFC. In particular, we have shown that the MA-DDPG framework allows for centralized learning with decentralized execution at different levels of observability and time requirement, which is a practical approach for real-world cellular environments. Furthermore, we have

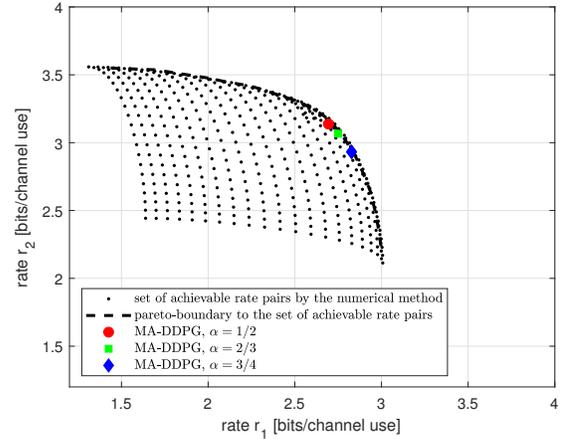
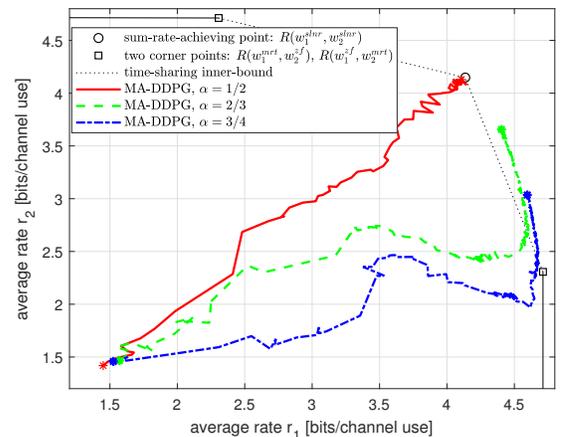
Fig. 3. Achieved rate pairs by MA-DDPG with  $\alpha = 1/2, 2/3,$  and  $3/4$ 

Fig. 4. Learning behaviors of MA-DDPGs in terms of the average rate pairs

addressed the phase ambiguity issue with the conventional baseband signal representation used in radio communications and proposed the phase ambiguity elimination method. The numerical simulation results show that the phase-ambiguity elimination in state space is crucial for successful training of MADRL in wireless communication systems. The proposed method can be also applied in precoding action space.

## REFERENCES

- [1] G. Wikström, J. Peisa, P. Rugeland, N. Johansson, S. Parkvall, M. A. Girnyk, G. Mildh, and I. L. da Silva, "Challenges and Technologies for 6G," in *6G Wireless Summit*, 2020.
- [2] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," in *International Conference on Learning Representations (ICLR)*, 2016.
- [3] R. Lowe, Y. Wu, Pritzel, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS)*, pp. 6382–6393, December 2017.
- [4] J. N. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," in *AAAI Conference on Artificial Intelligence (AAAI)*, pp. 2094 – 2100, February 2018.

- [5] E. Dahlman, S. Parkval, and J. Skold, *4G LTE-Advanced Pro and The Road to 5G*. Third Edition, Academic Press, 2016.
- [6] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Second Edition, MIT Press, Cambridge, Massachusetts, London, 2017.
- [7] H. Lee, M. Girnyk, and J. Jeong, "Deep MIMO Autoprecoder," in *Proceedings of IEEE International Conference on Communications*, June 2020.
- [8] H. Lee, M. Girnyk, and J. Jeong, "Deep reinforcement learning approach to MIMO precoding problem: Optimality and Robustness," *submitted to IEEE Transactions on Wireless Communications*, June 2020.
- [9] E. Jorswieck, E. Larsson, and D. Danev, "Complete Characterization of the Pareto Boundary for the MISO Interference Channel," *IEEE Transactions on Signal Processing*, vol. 56, pp. 5292–5296, October 2008.
- [10] S.-H. Park, H. Park, and I. Lee, "Distributed Beamforming Techniques for Weighted Sum-Rate Maximization in MISO Interference Channels," *IEEE Communications Letters*, vol. 14, pp. 1131–1133, December 2010.
- [11] M. Sadek, A. Tarighat, and A. H. Sayed, "A Leakage-Based Precoding Scheme for Downlink Multi-User MIMO Channels," *IEEE Transactions on Wireless Communications*, vol. 6, no. 5, pp. 1711–1721, 2007.