



北京大學  
PEKING UNIVERSITY

# Meta-Critic Reinforcement Learning for IOS-Assisted Multi-User Communications in Dynamic Environments

Qinpei Luo<sup>\*</sup>, Boya Di<sup>\*</sup>, Zhu Han<sup>†</sup>

<sup>\*</sup> State Key Laboratory of Advanced Optical Communication Systems and Networks, School of Electronics, Peking University

<sup>†</sup> Electrical and Computer Engineering Department, University of Houston, TX, USA

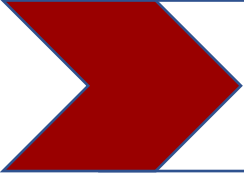
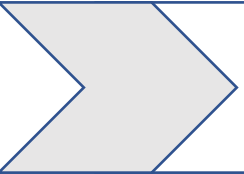
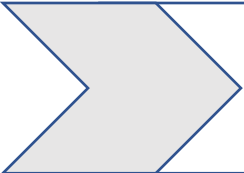
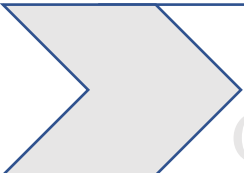
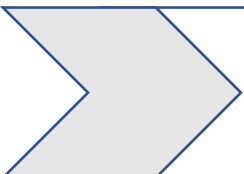


# Syllabus

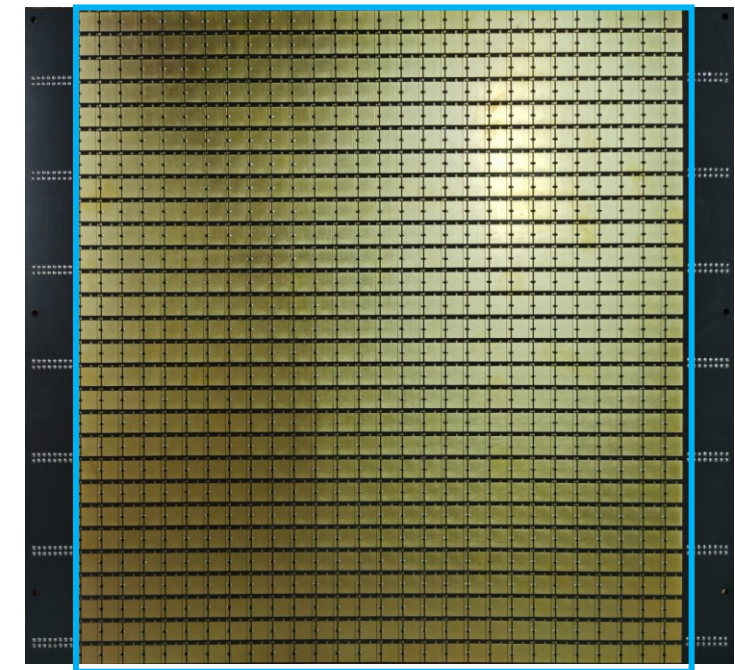
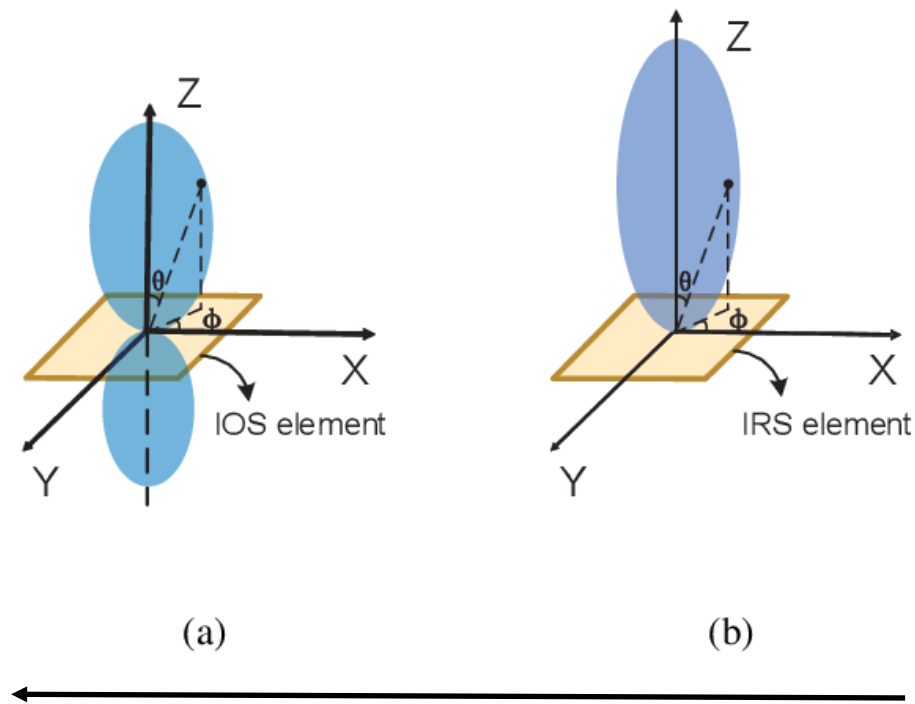
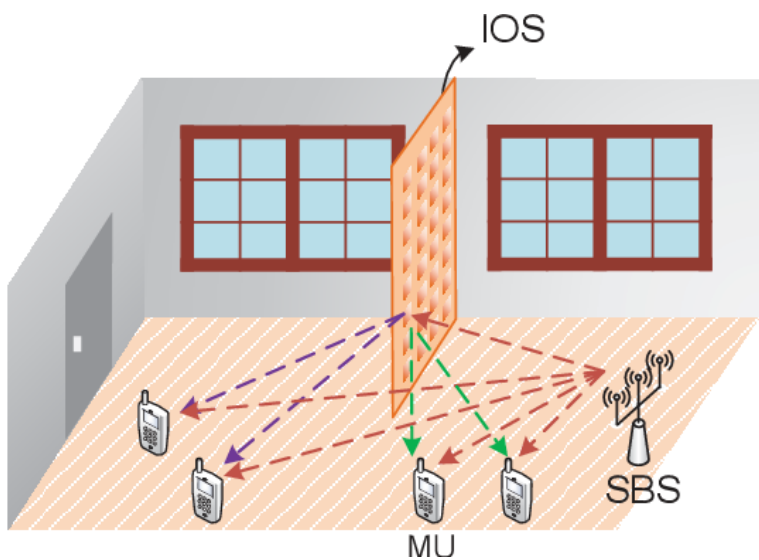
- Introduction to IOS
- Related Works and Limitations
- System Model & Problem Formulation
- MC-DDPG: A Faster Method for IOS Configurations in Dynamic Environment
- Simulation Results & Conclusion



# Agenda

-  Introduction to IOS
-  Related Works and Limitations
-  System Model & Problem Formulation
-  MC-DDPG: A Faster Method for IOS Configurations in Dynamic Environment
-  Simulation Results & Conclusion

# What is IOS? A promising solution to enhance the capacity of wireless networks



Intelligent Omni-Surface (IOS)  
Simultaneously **Reflection & Refraction**

Reflective Intelligent Surface (RIS)  
Only **Reflection** of incident signal

\*Source: Zhang, S., Zhang, H., Di, B., Tan, Y., Renzo, M.D., Han, Z., Poor, H.V., & Song, L. (2020). Intelligent Omni-Surface: Ubiquitous Wireless Transmission by Reflective-Transmissive Metasurface. *ArXiv, abs/2011.00765*.



# Challenges for Implementation of IOS

- **Numerous IOS elements**
  - Phase shifts of all of IOS elements need to be configured simultaneously, which brings difficulty in solution searching.
- **Dynamic Environment**
  - The channel state of environment changes rapidly, which requires real-time updates of IOS configuration.
- The above two things combines together to require an efficient beamforming scheme to tackle numerous IOS elements adapting to the varying channel information, users' positions, etc.



# Agenda

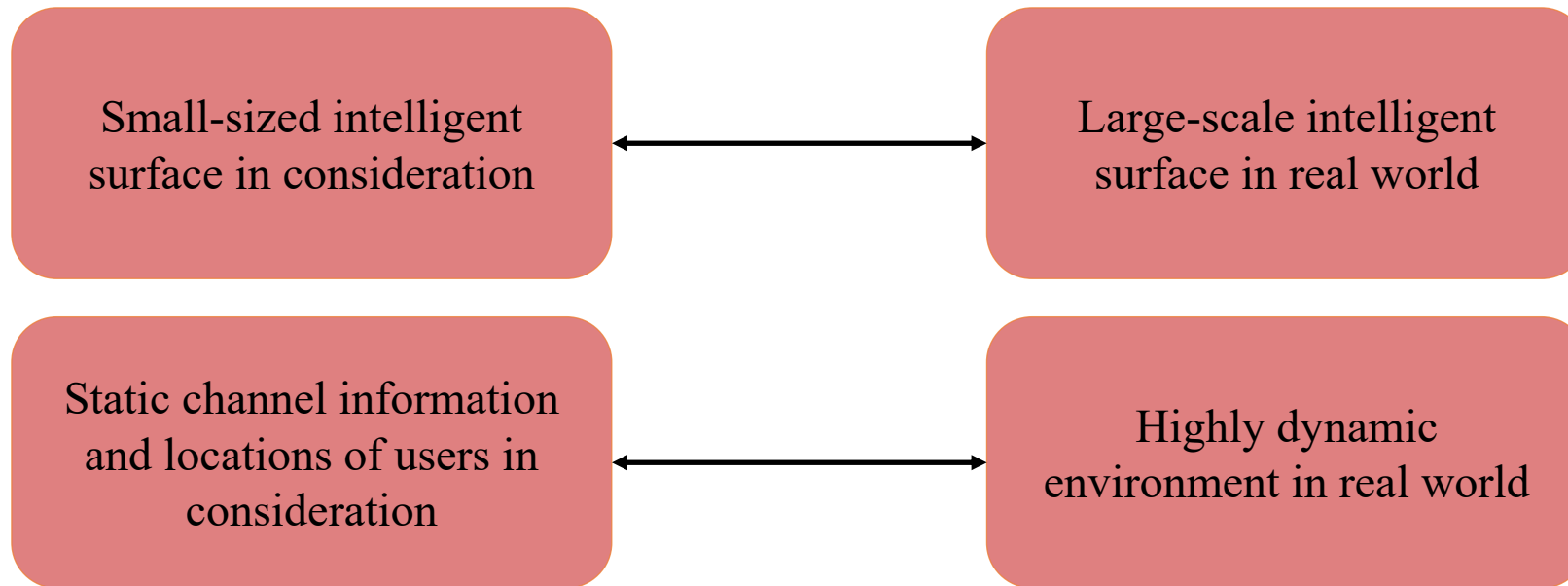
- Introduction to IOS
- Related Works and Limitations**
- System Model & Problem Formulation
- MC-DDPG: A Faster Method for IOS Configurations in Dynamic Environment
- Simulation Results & Conclusion



# Machine-Learning Based Beamforming

- Why ML is widely used?
  - Advanced ability in **extracting features from channel state information**.
- Reinforcement learning (RL) Method
  - Able to **well depict the interaction process** between intelligent surface and the environment.
  - HUANG, et al. (2020) develop a Deep RL based method to jointly design the transmit beamforming matrix and phase shifts of RIS.
  - LEE, et al. (2020) also use DRL to solve the problem of energy efficiency optimization.
  - ZHANG, et al. (2022) consider a system with multiple RISs and design a hierarchal policy network to improve the sum rate.

# Limitations of Current Methods



- We aim to develop an efficient beamforming scheme to address practical concerns
  - How to adapt to the dynamic case where the channel information and user positions vary with time?
  - How to deal with the numerous phase shift variables brought by a large-scale IOS in this case?

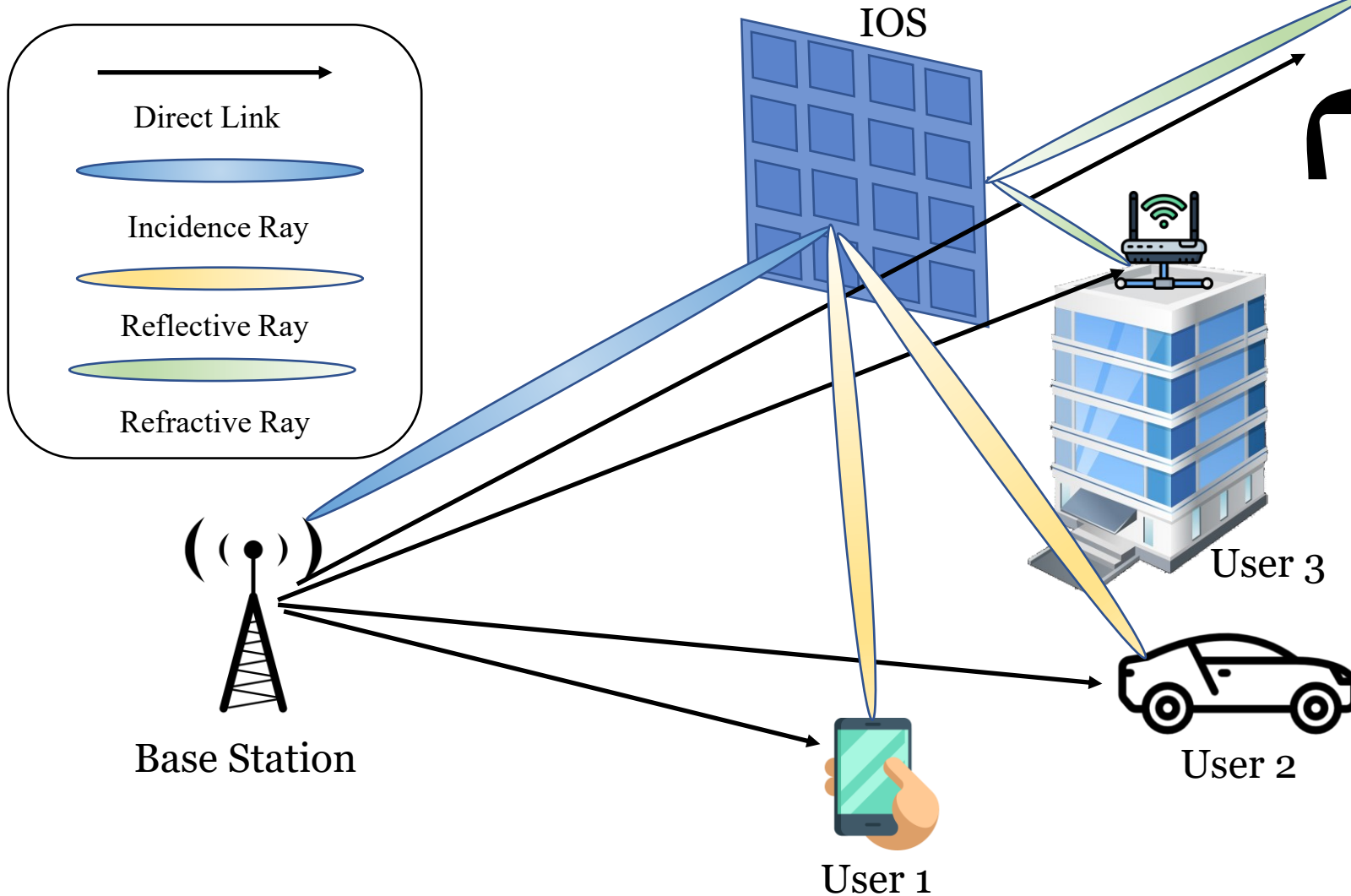




# Agenda

- Introduction to IOS
- Related Works and Limitations
- System Model & Problem Formulation**
- MC-DDPG: A Faster Method for IOS Configurations in Dynamic Environment
- Simulation Results & Conclusion

# Scenario Description



- MISO system:
  - $M$  antennas at Base Station (BS)
  - $K$  users with single antenna
- IOS:
  - Consisting of  $N$  elements
  - Able to reflect and refract the transmit signal
- Dynamic Environment
  - Channel states and locations of users may vary with time



# Channel Model

- For each user  $k$  we consider the Light-of-sight channel as a hybrid channel

$$\mathbf{H}_k^{LOS} = \Delta^u \mathbf{H}_{IU,k} \mathbf{\Theta} \mathbf{H}_{BI} + \mathbf{H}_{BU,k}$$

Where  $u \in \{r, t\}$  refers to the reflective and refractive respectively, while  $\Delta^u$  represents the energy split for each type of users.  $\mathbf{\Theta} = \text{diag}\{[e^{j\theta_1}, e^{j\theta_2}, \dots, e^{j\theta_N}]\}$  stands for the phase shifts of IOS.

- According to Saleh-Valenzuela Model, the channel of IOS-user, BS-IOS and BS-user can be further written into

$$\mathbf{H}_{BI} = \sqrt{S_1} \mathbf{A}_I \mathbf{\Sigma}_{BI} \mathbf{D}_B^H, \mathbf{H}_{IU} = \sqrt{S_{2,k}} \mathbf{A}_{IU,k} \mathbf{\Sigma}_{IU,k} \mathbf{D}_{I,k}^H, \mathbf{H}_{BU} = \sqrt{S_{3,k}} \mathbf{A}_{BU,k} \mathbf{\Sigma}_{BU,k} \mathbf{D}_{B,k}^H$$

In which  $\mathbf{A}$  and  $\mathbf{D}$  refers to transmit/receive steering matrices, the  $i$ -th column of each matrix is the steering vector and can be expressed by  $f(M, \theta) = \frac{1}{\sqrt{M}} [1, e^{j\pi\theta}, \dots, e^{j\pi(M-1)\theta}]^H$  where  $M$  is the number of antennas and  $\theta$  is the Angle-of-Arrival (AoA) or Angle-of-Departure (AoD).  $\mathbf{\Sigma}$  represents the gain of each channel, while  $S$  stands for the path loss.

- We assume the equivalent channel of each user follows Rician Distribution, i.e.,

$$\mathbf{H}_k = \sqrt{\frac{K^R}{1 + K^R}} \mathbf{H}_k^{LOS} + \sqrt{\frac{1}{1 + K^R}} \mathbf{H}_k^{NLOS}$$

$K^R$  is the Rician factor.  $\mathbf{H}_k^{NLOS}$  has similar expression as  $\mathbf{H}_k^{LOS}$ , but its AoDs or AoAs are randomly generated.



# Finite State Markov Channel

- We choose to fix the LOS component and discretize the NLOS component  $\mathbf{H}_k^{NLOS}$  into  $L$  levels.
- $\mathcal{H} = \{\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_k\}$
- Transition probability matrix:  $\mathbf{P} = \begin{pmatrix} p_{1,1} & \cdots & p_{1,L} \\ \vdots & \ddots & \vdots \\ p_{L,1} & \cdots & p_{L,L} \end{pmatrix}$
- $p_{l,l'} = Prob[\mathbf{H}_{t+1} = \mathbf{H}_{l'} | \mathbf{H}_t = \mathbf{H}_l], \mathbf{H}_l, \mathbf{H}_{l'} \in \mathcal{H}$
- $\mathbf{P}$  is generated randomly, so do the NLOS components.



# Sum Rate Maximization Formulation

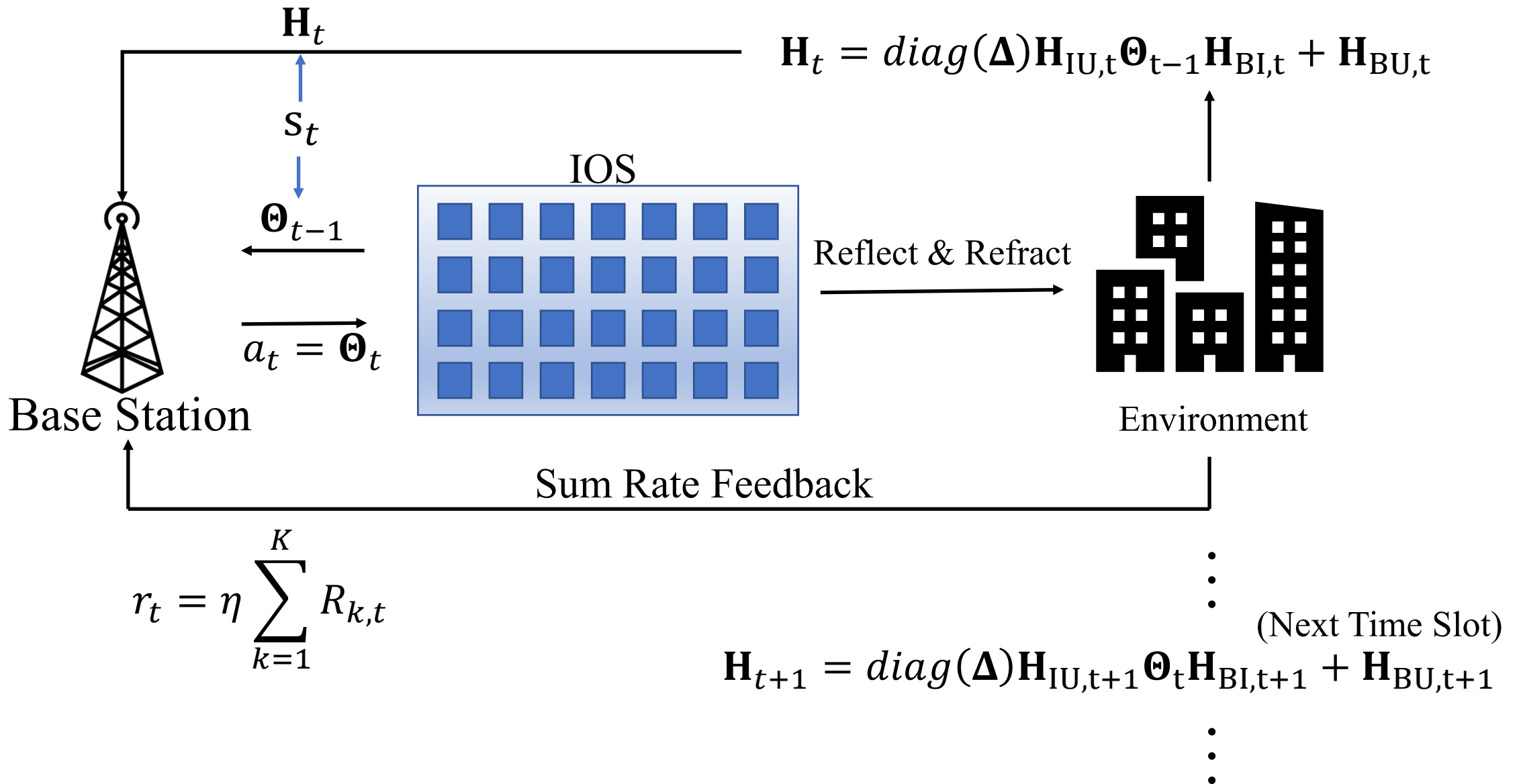
- We consider the sum rate maximization problem of all users in  $T$  time slots.

$$y_{k,t} = (\Delta_k \mathbf{H}_{IU,k} \Theta_t \mathbf{H}_{BI} + \mathbf{H}_{BU,k}) \sum_{j=1}^K \mathbf{V}_{j,t} m_j + n_{k,t}$$
$$\gamma_{k,t} = \frac{|(\Delta_k \mathbf{H}_{IU,k} \Theta_t \mathbf{H}_{BI} + \mathbf{H}_{BU,k}) \mathbf{V}_{k,t} m_k|^2}{(\Delta_k \mathbf{H}_{IU,k} \Theta_t \mathbf{H}_{BI} + \mathbf{H}_{BU,k}) \sum_{j=1, j \neq k}^K \mathbf{V}_{j,t} m_j + n_{k,t}}$$
$$R_{k,t} = \log_2(1 + \gamma_{k,t})$$
$$\mathbf{P1}: \max_{\{\mathbf{V}_t, \Theta_t\}} \sum_t \sum_k R_{k,t}$$

- Solving  $\mathbf{V}_t$  with fixed digital beamforming method as water-filling and zero-forcing, we can rewrite the problem as

$$\mathbf{P2}: \max_{\Theta_t} \sum_t \sum_k R_{k,t}$$

# Markov Decision Process Reformulation





# Agenda

- Introduction to IOS
- Related Works and Limitations
- System Model & Problem Formulation
- MC-DDPG: A Faster Method for IOS Configurations in Dynamic Environment**
- Simulation Results & Conclusion

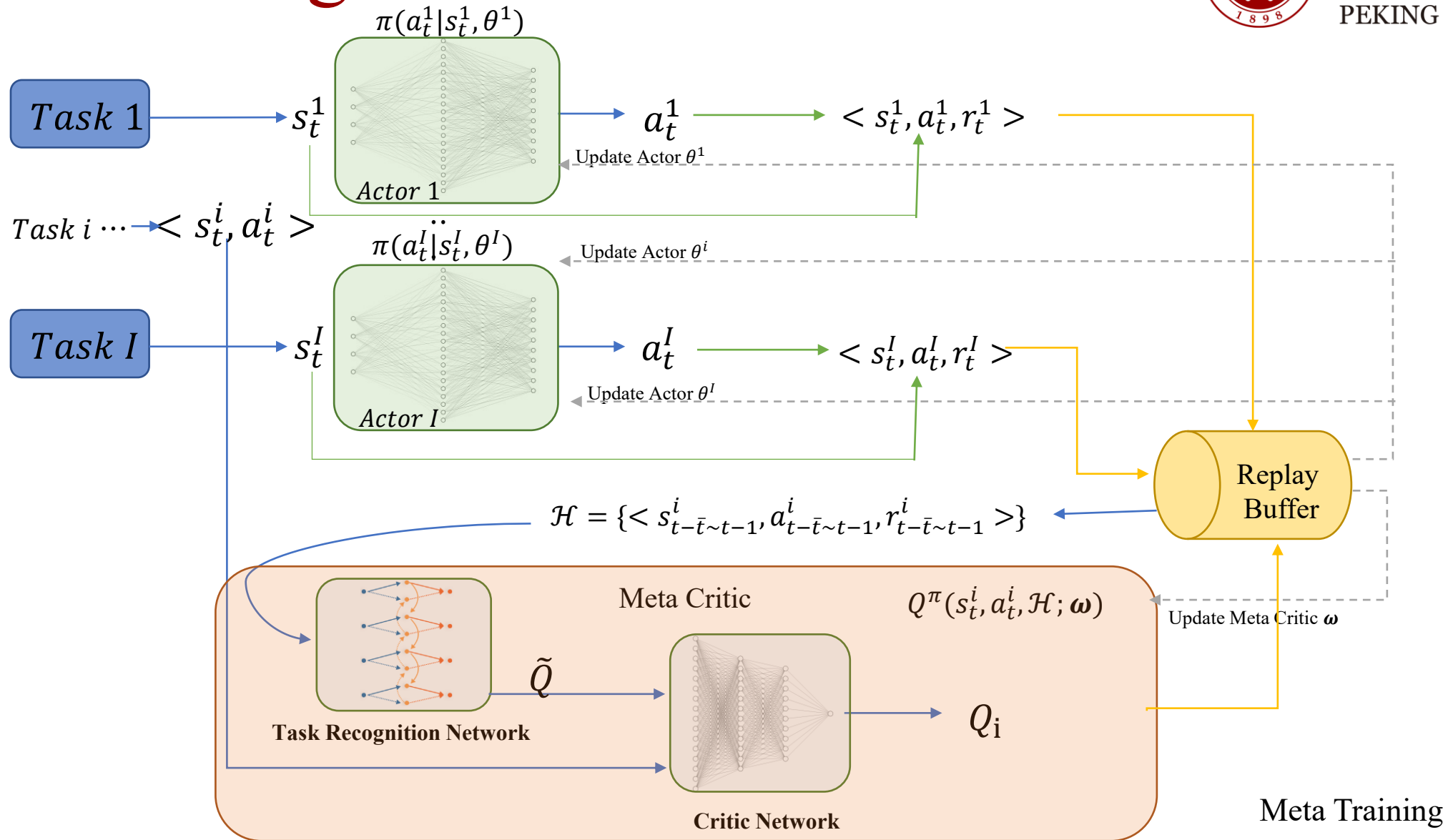


# Basic Definition

- **Task:** Denotes a process of the BS maximizing the sum rates of all users in a fixed number of time slots. For different tasks, the parameters of BS and IOS are set as the same, while the channel states and locations of users are various.
- **Actor:** It receives the information of state from the task in each time slot and outputs correspondent action. We adopt a neural network as the policy of the actor.
- **Meta-Critic:** Consists of two parts, a task recognition network and a critic network. The former extracts the history information and generates the task-recognition Q-value, while the critic network outputs a task-specific Q-value to update the actor networks.

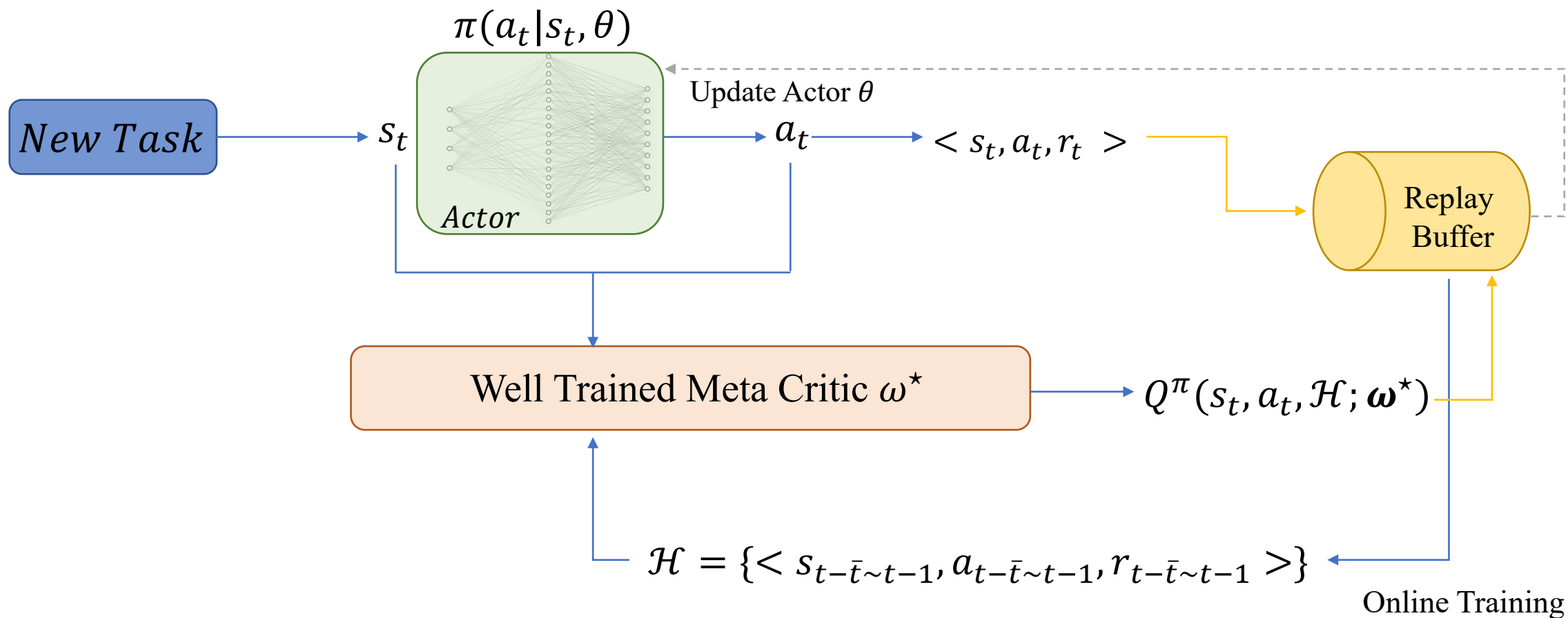


# Meta Learning Phase



In the meta-learning phase, both the actors and meta-critic are updated.

# Online Learning Phase



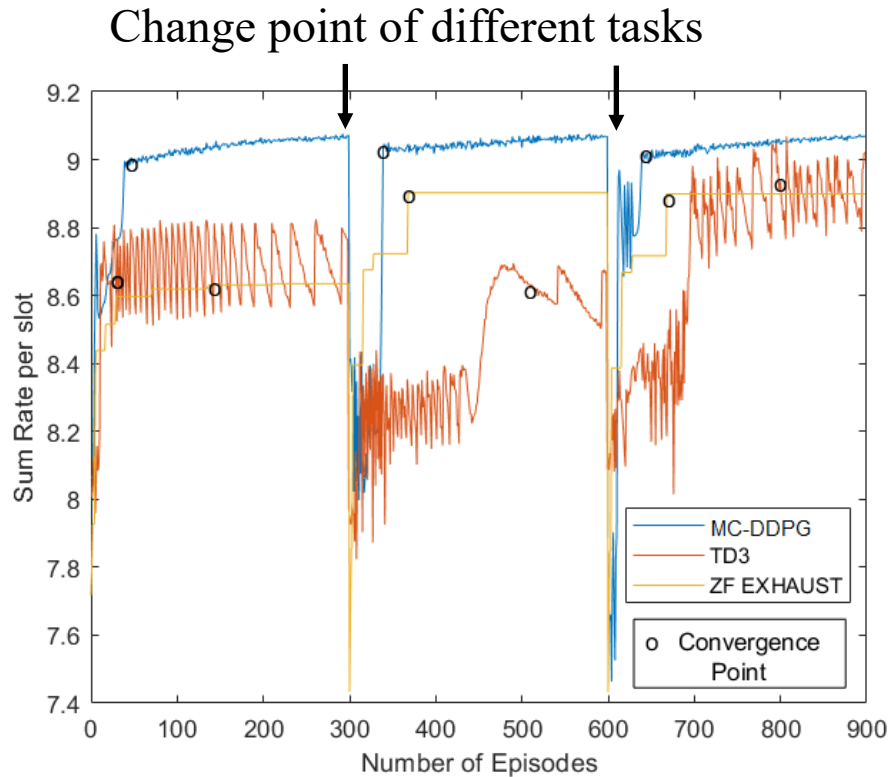
In the online learning phase, only the actors are updated, while the well-trained meta-critic is kept static.



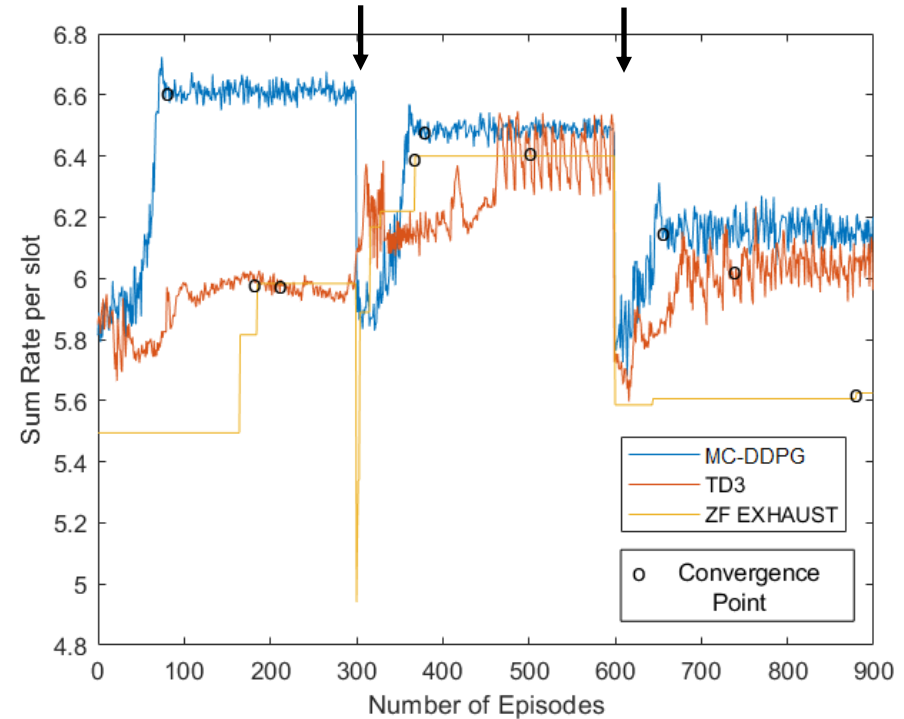
# Agenda

- Introduction to IOS
- Related Works and Limitations
- System Model & Problem Formulation
- MC-DDPG: A Faster Method for IOS Configurations in Dynamic Environment
- Simulation Results & Conclusion**

# Sum Rate Performance in Dynamic Settings

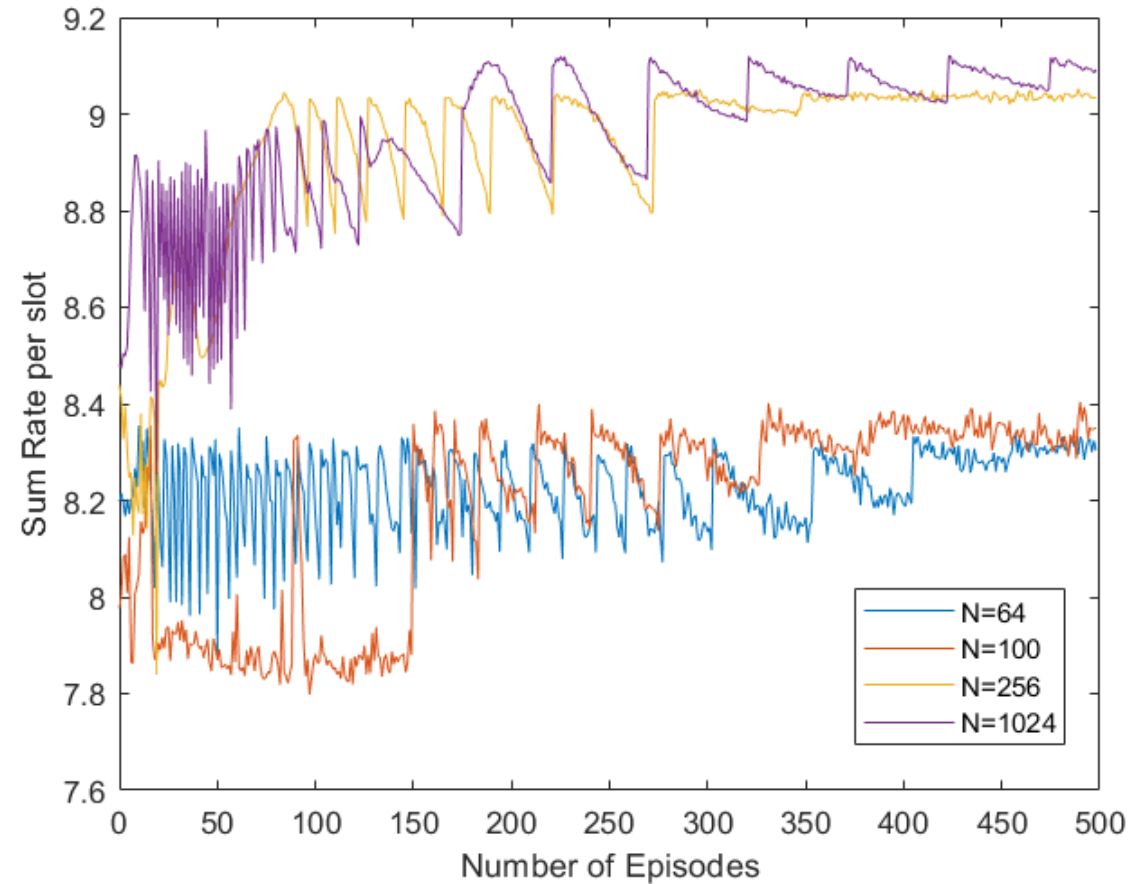


Sum Rate Performance with respect to the varying channel states



Sum Rate Performance with respect to User's locations

# Influence of the Number of IOS Elements





# Conclusion

- Current works seldom consider the challenges brought by a large number of IOS elements and dynamic environment.
- We proposed MC-DDPG, a meta-critic RL scheme for sum rate maximization given the limited CSI, which is able to:
  - Achieve a **faster convergence speed** and a **higher sum rate** compared to the benchmarks.
  - The **robustness of MC-DDPG against IOS sizes** is verified.
- We can draw two take-away conclusions:
  - The designed meta-critic significantly **enhances the robustness** of the IOS-assisted multi-user communications **against user mobility and the dynamic CSI**.
  - There exists a **trade-off between the convergence speed and the achievable sum rate** of MC-DDPG.



北京大學  
PEKING UNIVERSITY

# Thank you! Q&A

Qinpei Luo<sup>\*</sup>, Boya Di<sup>\*</sup>, Zhu Han<sup>†</sup>

<sup>\*</sup> State Key Laboratory of Advanced Optical Communication Systems and Networks, School of Electronics, Peking University

<sup>†</sup> Electrical and Computer Engineering Department, University of Houston, TX, USA