# DRL-Aided Joint Resource Block and Beamforming Management for Cellular-Connected UAVs

**Dr. Yuanjian Li**, Prof. Mathini Sellathurai, Dr. Zheng Chu, Prof Pei Xiao, and Prof. Hamid Aghvami *Fellow, IEEE*

*Research Fellow* @ Nanyang Technological University (NTU), Singapore

December 6, 2023

NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

HERIOT WATT UNIVERSITY | 200 Years

UNIVERSITY OF SURREY

KING'S College LONDON

## Outline

## Brief Overview

In this paper, we investigate a cellular-connected unmanned aerial vehicle (UAV) network, where multiple UAVs receive messages from base stations (BSs) in the down-link, and in the meantime, BSs serve their paired ground user equipments (UEs). To effectively manage inter-cell interferences (ICIs) among UEs due to intense reuse of time-frequency resource block (RB) resource, a first $p$-tier based RB coordination criterion is adopted. Then, to enhance wireless transmission quality for UAVs while protecting terrestrial UEs from being interfered by ground-to-air (G2A) transmissions, a radio resource management (RRM) problem of joint dynamic RB coordination and time-varying beamforming design is formulated to minimize UAV's ergodic outage duration (EOD).

## Brief Overview

To cope with conventional optimization techniques' inefficiency in solving the formulated RRM problem, a deep reinforcement learning (DRL)-aided solution is proposed, where deep double duelling Q network (D3QN) and twin delayed deep deterministic policy gradient (TD3) are invoked to deal with RB coordination in the discrete action domain and beamforming design in the continuous action regime, respectively. Numerical results illustrate the effectiveness of the proposed hybrid D3QN-TD3 algorithm, compared to representative baselines.

## Main Contribution

Radio resource management issue on interference coordination and beamforming design in down-link cellular-connected UAV networks is investigated in this paper. It is worth noting that the terrestrial transmissions between GUEs and BSs are protected to be not contaminated by the down-link G2A channels. The main contributions of this paper can be concluded as follows.

- A joint time-frequency resource block (RB) allocation and beamforming design optimization problem is formulated to minimize the ergodic outage duration (EOD) of UAV, for arbitrary given trajectory. Specifically, the RB allocation is utilized to assign proper RB resource to UAVs while insuring that the terrestrial transmissions are not violated by the potential co-channel interferences generated from BSs aligned to serve UAVs. To enhance the quality of received signals at UAVs after RB allocation, transmit beamforming design is invoked in the presence of imperfect G2A channel estimation.

## Main Contribution

- To deal with the difficulty of traditional optimization-based methods solving the proposed EOD minimization problem, a deep reinforcement learning (DRL) aided solution is proposed via mapping the proposed EOD minimization problem into an outer Markov decision process (MDP) and an inner MDP. The outer MDP reflects the dynamic RB possession environment at BSs, while the inner MDP tracks the corresponding small-scale fading characteristics. The outer MDP contains discrete action space (i.e., RB indices), which is tackled by invoking deep double duelling Q network (D3QN), while the continuous action space (i.e., beamforming vectors) in the inner MDP is dealt with twin delayed deep deterministic policy gradient (TD3) approach. The proposed hybrid D3QN-TD3 algorithm can optimize EOD performance for UAVs via interactively interacting with the outer and inner environments, of which the D3QN and TD3 agents are able to offer independent EOD performance gains.

## Main Contribution

- In contrast to the majority of related literature adopting statistical G2A channel model (e.g., probabilistic G2A channel model), LoS/NLoS G2A pathloss is determined via checking potential blockages between UAV and BSs in this paper, according to one realization of local building distribution suggested by the International Telecommunication Union (ITU). The considered G2A channel model is more practical than its statistical counterpart which can only reflect average pathloss gain over large number of similar building distribution realizations because building distribution in local area remains unchanged in practice.
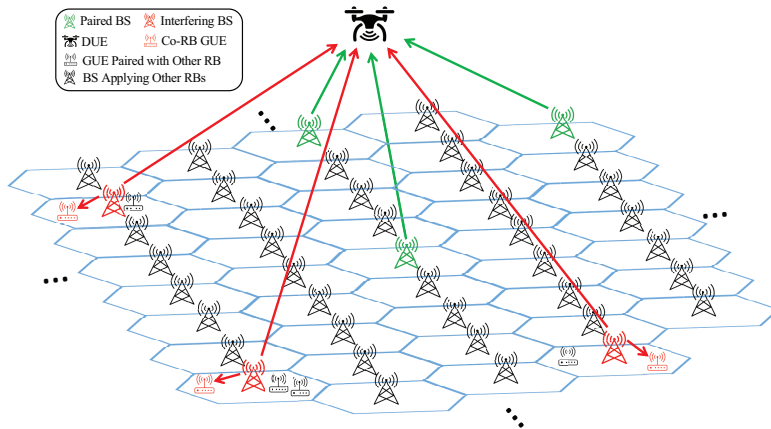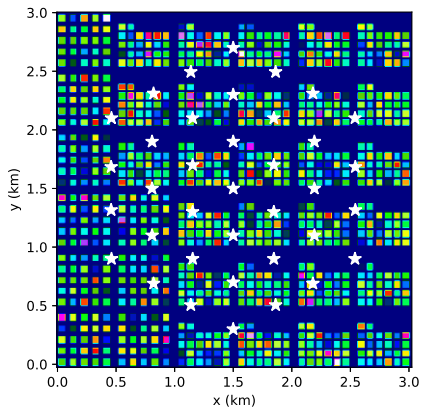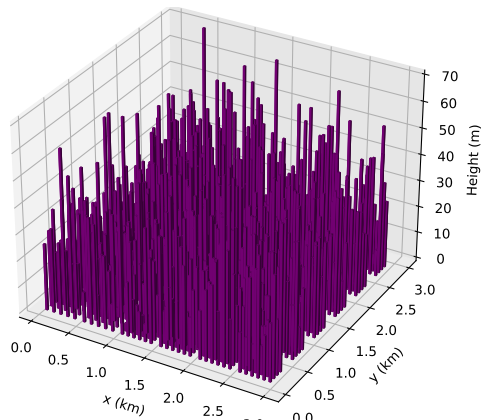
Brief Overview
oo

Main Contribution
ooo

System Model
●ooo

Problem Formulation
ooo

The Proposed DRL Algorithm
oooooo

Simulation Results
ooo

## System Model



Figure 2: System model

Brief Overview
○○

Main Contribution
○○○

System Model
○●○○

Problem Formulation
○○○

The Proposed DRL Algorithm
○○○○○○

Simulation Results
○○○

# BS Distribution



(a) First $p$-tier BSs

(b) An example of BS grouping

Figure 3: Illustrations of system model, the first $p$-tier set of BS $b$ and an instance of BS grouping in the case of $p = 1$, where $card(\mathcal{TI}_b(1)) = 7$, $card(\mathcal{TI}_b(2)) = 19$ and $card(\mathcal{TI}_b(3)) = 37$

# Building Distribution



(a) 2D distribution of local buildings and BSs

(b) 3D view of local building distribution

Figure 4: The considered building distribution

## Channel Model

$$\text{PL}_l = \begin{cases} 28.0 + 22\log_{10}(d_{bu}) + 20\log_{10}(f_c), & l = \text{LoS} \\ -17.5 + [46 - 7\log_{10}(h)] \times \log_{10}(d_{bu}) + 20\log_{10}\left(\frac{40\pi f_c}{3}\right), & l = \text{NLoS} \end{cases}, \quad (1)$$

in which $r_{bu} = \sqrt{d_{bu}^2 - h^2}$, $\varepsilon_1 = \max\{460\log_{10}(h) - 700, 18\}$, $\varepsilon_2 = 4300\log_{10}(h) - 3800$, $f_c$ represents the carrier frequency and $d_{bu} = ||\vec{q}_u - \vec{q}_b||_2$ calculates the Euclidean distance between DUE $u$ and ground BS $b$. Since the proposed design on beamforming vectors aims to be adaptive to arbitrary small-scale fading environment, we denote $\vec{h}_{bu} \in \mathbb{C}^{1 \times M}, \forall b \in \mathcal{B}, u \in \mathcal{U}$ as the small-scale fading component for B2D channels and pose no assumptions on its modelling, while an example of specific B2D small-scale fading model will be discussed in numerical result section.

## The Received Signal

The received signal of DUE $u$ over RB $k$ at time $t$ can be given by

$$y_u^k(t) = C_u^k(t) \left[ \sum_{b \in \breve{\mathcal{B}}_o^k} \sqrt{10^{\frac{-\mathsf{PL}_l}{10}}} \, \vec{h}_{bu} \vec{w}_{bu} x_u(t) + \sum_{b \in \mathcal{B}_o^k} \sqrt{10^{\frac{-\mathsf{PL}_l}{10}}} \, \vec{h}_{bu} \vec{w}_{bg} x_{bg}(t) + n_u^k \right], \quad (2)$$

where $\vec{w}_{bu} \in \mathbb{C}^{M \times 1}$ indicates the transmit beamforming vector at BS $b \in \breve{\mathcal{B}}_o^k$ for DUE $u$, $\vec{w}_{bg} \in \mathbb{C}^{M \times 1}$ represents the transmit beamforming vector at BS $b \in \mathcal{B}_o^k$ for corresponding GUEs, $x_u(t) \sim \mathcal{CN}(0, P)$ is the intended message from BS $b$ to DUE $u$, $x_{bg}(t) \sim \mathcal{CN}(0, P)$ implies the signal for GUEs, and $n_u^k \sim \mathcal{CN}(0, \sigma^2)$ denotes the received additive complex Gaussian noise (AWGN) at DUE $u$.

## SINR

Denote $C_u^k(n) \in \{0, 1\}$ as the RB association indicator which means that DUE $u$ is occupying RB $k$ over time slot $n$ when $C_u^k(n) = 1$, and $C_u^k(n) = 0$ otherwise. Each DUE is assumed to occupy at most one single RB each time, then we have $\sum_{k=1}^{K} C_u^k(n) \leq 1$. Then, instantaneous SINR of DUE $u$ at time $t$ can be calculated as

$$\Gamma_u(t) = \sum_{k=1}^{K} \frac{C_u^k(t) \left[ \sum_{b \in \breve{\mathcal{B}}_o^k} \sqrt{P 10^{\frac{-\text{PL}_l}{10}}} |\vec{h}_{bu} \vec{w}_{bu}| \right]^2}{I_u^k(t) + \sigma^2}, \tag{3}$$

where $I_u^k(t) = \sum_{b \in \mathcal{B}_o^k} P 10^{\frac{-\text{PL}_l}{10}} |\vec{h}_{bu} \vec{w}_{bg}|^2$ means ICIs introduced by co-channel BSs from $\mathcal{B}_o^k$.

This paper focuses on minimizing ergodic outage duration (EOD) suffered from DUE's flight, approximately given by

$$(\text{P1}): \min_{\substack{C_u^k(n), \\ \vec{w}_{bu}(n,i)}} \sum_{n=1}^{N_u} \sum_{i=1}^{\varsigma} \frac{\delta_u}{\varsigma} ITOP_u\{C_u^k(n), \vec{w}_{bu}(n,i); \hat{h}(n,i)\}, \tag{4a}$$

$$\text{s.t.} \Sigma_{k=1}^{K} C_u^k(n) \leq 1, \forall n \in [1, N_u], \tag{4b}$$

$$C_u^k(n) \in \{0, 1\}, \forall k \in \mathcal{K}, \forall n \in [1, N_u]. \tag{4c}$$

$$||\vec{w}_{bu}(n,i)||^2 = 1, \forall b \in \breve{\mathcal{B}}_o^k, \forall n \in [1, N_u], \tag{4d}$$

where constraint (4d) ensures that the transmit power of each available BS $b \in \breve{\mathcal{B}}_o^k$ equals $P$. Denote $\Gamma_{th}$ as the SINR threshold, indicator of transmission outage probability (TOP) is defined as $ITOP_u\{C_u^k(n), \vec{w}_{bu}(n,i); \hat{h}(n,i)\} = 1$ when $\Gamma_u(n,i) < \Gamma_{th}$, and $ITOP_u\{C_u^k(n), \vec{w}_{bu}(n,i); \hat{h}(n,i)\} = 0$ otherwise, where $\hat{h}(n,i)$ and $\vec{w}_{bu}(n,i)$ indicate one realization of fast fadings and corresponding beamforming vectors in time slot $n$, respectively. Then, to calculate the corresponding TOP, $\varsigma \gg 1$ times of SINR measurement are performed. The length of time slot is defined as $\delta_u$ and thus the number of time slots equals $N_u = T_u/\delta_u$ for the DUE $u$.

The Proposed DRL Solution

To derive a flexible solution solving (4) in a dynamic RB possession (RBP) and time-varying small-scale fading scenario, a deep reinforcement learning (DRL)-aided solution is proposed after mapping the considered optimization problem into two Markov decision processes (MDPs). Specifically, the deep double duelling netwrok (D3QN) maps the outer state and the RB selection into Q values, while actor of twin delayed deep deterministic policy gradient network (TD3) agent transforms the inner state to beamforming vector and TD3's critic network evaluates the corresponding Q values.
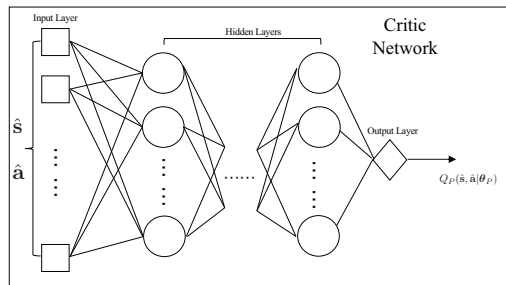
## CNN-Attached Duelling DQN



Figure 5: Architecture of CNN-attached duelling DQN

## Architectures of TD3's Actor and Critic



(a) Layout of the actor network

(b) Layout of the critic network

Figure 6: Architecture of TD3 network

## Pseudocode of the Proposed Hybrid D3QN-TD3 Solution

---

**Algorithm 1:** The hybrid D3QN-TD3 solution

1. **Initialization:** Initialize D3QN $Q_{D3}(s, a|\boldsymbol{\theta}_{D3})$ and its target network $Q_{D3}(s, a|\boldsymbol{\theta}_{D3}^-)$, with $\boldsymbol{\theta}_{D3}^- \leftarrow \boldsymbol{\theta}_{D3}$. Initialize TD3, including actor network $\mu(s|\boldsymbol{\theta}_\mu)$, twin critics $Q_P(s, a|\boldsymbol{\theta}_{P_j})$, target actor $\mu(s|\boldsymbol{\theta}_\mu^-)$ and twin target critics $Q_P(s, a|\boldsymbol{\theta}_{P_j}^-)$, with $\boldsymbol{\theta}_\mu^- \leftarrow \boldsymbol{\theta}_\mu$ and $\boldsymbol{\theta}_{P_j}^- \leftarrow \boldsymbol{\theta}_{P_j}$. Initialize D3QN and TD3 replay buffers R and Ř with capacity Ď and Ĺ, respectively;

2. **for** $episode = [1, epi]$ **do**
3.     Initialize the outer environment and reset the UAV's location to $\vec{q}_u(0)$;
4.     **for** $i = [1, epo_{outer}]$ **do**
5.        Observe the outer state $\mathbf{s}_i$;
6.        Select the outer action $\mathbf{a}_i$, observe the available set $\breve{\mathcal{B}}_o^{\mathbf{a}_i}$ and the occupied set $\mathcal{B}_o^{\mathbf{a}_i}$;
7.        Randomly select a BS $\check{b} \in \breve{\mathcal{B}}_o^{\mathbf{a}_i}$ and check its B2D pathloss type, i.e., LoS or NLoS, then initialize the inner environment;
8.        **for** $j = [1, epo_{inner}]$ **do**
9.           Observe the inner state $\hat{\mathbf{s}}_j$;
10.           Select and execute the inner action $\hat{\mathbf{a}}_j$, observe the next inner state $\hat{\mathbf{s}}_{j+1}$ and calculate the inner reward $\hat{\mathbf{r}}_j$;
11.           Store transition $(\hat{\mathbf{s}}_j, \hat{\mathbf{a}}_j, \hat{\mathbf{s}}_{j+1}, \hat{\mathbf{r}}_j)$ into Ř;
12.           **if** $card\,(\check{R}) \geq N_P$ **then**
13.              Sample a mini-batch of $N_P$ transitions from Ř, then update the twin critic networks $Q_P(s, a|\boldsymbol{\theta}_{P_j})$ via gradient descent method in (21);
14.              Every $N_{pud}$ times twin critics are trained, update actor network $\mu(s|\boldsymbol{\theta}_\mu)$ via gradient ascent in (26), and update target networks $Q_P(s, a|\boldsymbol{\theta}_{P_j}^-)$ and $\mu(s|\boldsymbol{\theta}_\mu^-)$, following Polyak averaging rule in (27) and (28), respectively;
15.     Execute the outer action $\mathbf{a}_i$, observe the next outer state $\mathbf{s}_{i+1}$ and calculate the outer reward $\mathbf{r}_i$;
16.     Store transition $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}_{i+1}, \mathbf{r}_i)$ into R;
17.     **if** $card\,(R) \geq N_{D3}$ **then**
18.        Sample a mini-batch of $N_{D3}$ transitions from R, update D3QN network $Q_{D3}(s, a|\boldsymbol{\theta}_{D3})$ via gradient descent in (16);
19.        Update the D3QN target network $Q_{D3}(s, a|\boldsymbol{\theta}_{D3}^-)$ every $\Upsilon_{D3}$ steps, i.e., $\boldsymbol{\theta}_{D3}^- \leftarrow \boldsymbol{\theta}_{D3}$;
20.     Update $\epsilon \leftarrow \epsilon \times dec_\epsilon$ and $\sigma_P^2 \leftarrow \sigma_P^2 \times dec_\sigma$;

---

Figure 7: Pseudocode of the proposed hybrid D3QN-TD3 solution

## Workflow Diagram of the Proposed Hybrid D3QN-TD3 Solution
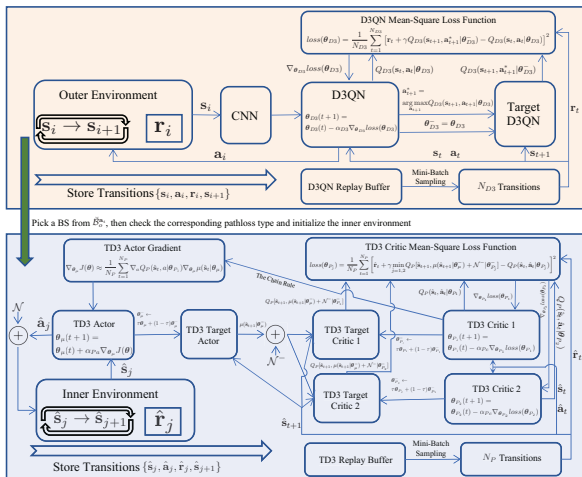


Figure 8: Workflow diagram

## Example of Offline Exploitation

Note that it is conducted on a central coordinator (C2) that controls BSs via high-speed fronthaul links and is genuinely an online-centralized-learning-and-offline-decentralized-execution algorithm, for realizing its efficient implementation without introducing heavy burden of computations or excessive delays or overheads of information transfer, e.g., imperfect CSIs and designed beamforming vectors between the C2 and available BSs, during its exploitation for RRM.
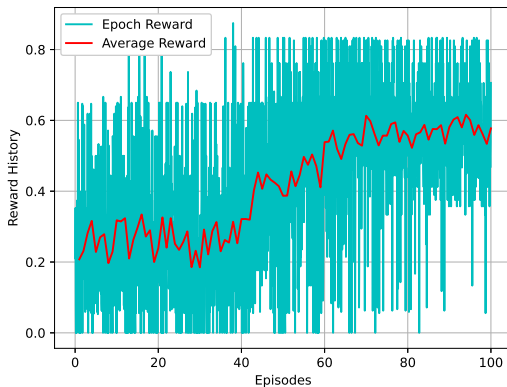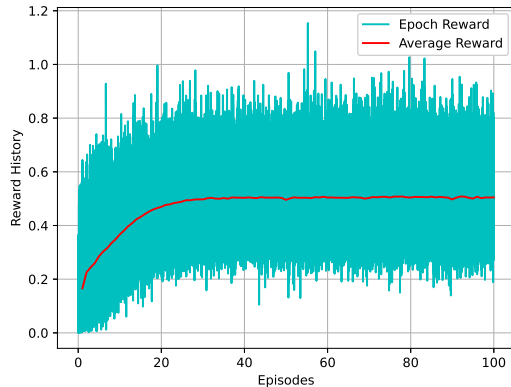


Figure 9: Offline exploitation
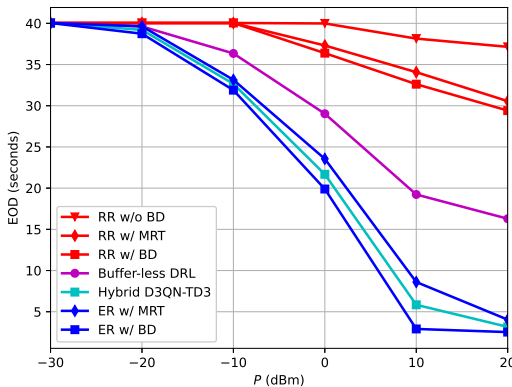
## Simulation Results
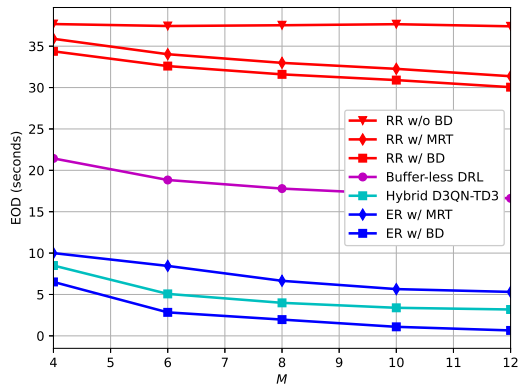


(a) Reward history of D3QN

(b) Reward history of TD3

Figure 10: Reward history

## Simulation Results



(a) Performance comparison vs $P$

(b) Performance comparison vs $M$

Figure 11: Performance comparison

Brief Overview
○○

Main Contribution
○○○

System Model
○○○○

Problem Formulation
○○○

The Proposed DRL Algorithm
○○○○○○

Simulation Results
○○●

**The End**

# Thanks for your attentions

This is the end of today's demonstration