# Deep Reinforcement Learning Optimized Intelligent Resource Allocation in Active RIS-Integrated TN-NTN Networks

Muhammad Ahmed Mohsin[1], Hassan Rizwan[2], Muhammad Jazib[3,4], Muhammad Iqbal[4], Muhammad Bilal[5], Tabinda Ashraf[4], Muhammad Farhan Khan[6], Jen-Yi Pan[4]

[1]Dept. of Electrical Engineering, Stanford University, USA
[2]Dept. of Electrical Engineering, University of California, Riverside, USA
[3]Dept. of Electrical Engineering, Pakistan Institute of Engineering and Applied Sciences, Pakistan
[4]Dept. of Communication Engineering, National Chung Cheng University, Taiwan
[5]Dept. of Electrical Engineering, University of California, Irvine, USA
[6]School of Computer Science and Information Technology, University College Cork, Ireland
Corresponding author emails: {muahmed@stanford.edu, hrizw002@ucr.edu, iqbal.marjan@gmail.com}

*Abstract*—**This work explores the deployment of active reconfigurable intelligent surfaces (A-RIS) in integrated terrestrial and non-terrestrial networks (TN-NTN) while utilizing coordinated multipoint non-orthogonal multiple access (CoMP-NOMA). Our system model incorporates a UAV-assisted RIS in coordination with a terrestrial RIS which aims for signal enhancement. We aim to maximize the sum rate for all users in the network using a custom hybrid proximal policy optimization (H-PPO) algorithm by optimizing the UAV trajectory, base station (BS) power allocation factors, active RIS amplification factor, and phase shift matrix. We integrate edge users into NOMA pairs to achieve diversity gain, further enhancing the overall experience for edge users. Exhaustive comparisons are made with passive RIS-assisted networks to demonstrate the superior efficacy of active RIS in terms of energy efficiency, outage probability, and network sum rate.**

*Index Terms*—**Active RIS, DRL, NOMA, Resource Allocation, TN-NTN**

## I. INTRODUCTION

The proliferation of wireless devices has led to a surge in demand for communication resources. Traditional orthogonal multiple access (OMA) schemes, which allocate dedicated time-frequency resources to each user, are becoming increasingly inefficient. NOMA has emerged as a promising technology that efficiently tackles the problem by leveraging superposition coding and successive interference cancellation (SIC) to enable multiple users to share the same time-frequency resources. This leads to enhanced spectral efficiency of a system [1], [2]. By allocating power to users based on channel conditions, NOMA can improve fairness and accommodate more users within a given spectrum.

In a multi-cell environment, cell-edge users often suffer from severe inter-cell interference as they are located closer to the boundaries of their cells and are more susceptible to interference from neighboring cells. To mitigate this issue, CoMP transmission is utilized to enable coordinated joint detection and joint transmission among multiple base stations, which can significantly improve the performance of cell-edge users [3].

Terrestrial RIS is deployed in densely populated urban areas where a stable user demand persists, while A-RIS serves dynamic environments where user demands fluctuate rapidly [4], [5]. Moreover, in an urban environment, there will be obstacles for the UAV that create abandoned areas in the model which are a no-fly zone for UAV. Therefore, it is imperative to optimize UAV trajectory to improve the sum rate and ensure UAV safety. Other works utilize the UAV as a BS which offers lesser flexibility than an RIS [6], since RISs can be dynamically reconfigured to adjust their reflecting properties based on changing channel conditions and even direct beams to users for better network capacity.

Typically, UAV communication systems use passive RISs and optimize their phase shift matrix to improve system throughput. However, utilizing an active RIS and optimizing its amplification factor matrix will substantially impact the overall sum rate of the environment as the comparison in [7] concludes that active RIS overcomes the multiplicative fading effect found in passive RIS to ultimately, provide a better sum rate gain.

While DRL has been increasingly applied to optimize systems involving RISs and UAVs, existing research primarily focuses on either continuous or discrete action spaces [8]. To achieve optimal performance, a hybrid approach that combines both continuous and discrete action spaces is essential. DRL algorithms that can use both action spaces simultaneously would offer promising solutions for optimization. Furthermore, the optimization of the active RIS amplification matrix using DRL has not yet been explored.

Keeping in view the literature gap as motivation, the contributions of this paper are threefold. A hybrid DRL agent, H-PPO, is utilized to optimize phase shifts, UAV trajectory, NOMA power allocation factors and base station transmit power to achieve maximum sum rate under optimal SNRs.

Phase shifts for terrestrial RIS are optimized to cancel the interference from the non-CoMP base station at the edge user to maximize user experience and network capacity at the edge. Exhaustive results comparisons are made with passive RIS to demonstrate the efficacy of A-RIS in different optimization settings for various performance metrics like outage, sum rate, and fairness.

## II. SYSTEM MODEL & PROBLEM FORMULATION

### A. System Description

We consider an active RIS-assisted CoMP-NOMA downlink SISO transmission network with UAV deployment. The network is distributed in $M$ circular grids with a single antenna BS at the center of each grid $m$, where $m \in \mathcal{M} = \{1, 2, ..., M\}$. Every $\text{BS}_m$ in grid $m$ utilizes two-user NOMA to provide a downlink channel to center users and edge users. The center users reside within the radius of the grid while being served by $\text{BS}_m$ and are denoted as $\text{U}_c^m$, where $c \in \mathcal{C} = \{1, 2, ..., C\}$ represents all center users. On the other hand, the edge users reside outside the grid radius and are denoted as $\text{U}_e^m$ served by $\text{BS}_m$, where $e \in \mathcal{E} = \{1, 2, ..., E\}$ represents all edge users. All users in the system can be denoted as $\mathcal{U} = \mathcal{C} \cup \mathcal{E}$, while each user $u \in \mathcal{U}$ is treated as either an edge user or a center user. The total UAV flight time is divided into $t$ time slots where $t \in \mathcal{T} = \{1, ..., t_s\}$, where $t_s$ is the total flight time for UAV. Moreover, UAV cannot fly in the forbidden zones present in the system model which are modeled as circular disks of radius $d_{\min}$ centered around different obstacles where an obstacle $o \in \mathcal{O} = \{1, 2, ..., O\}$.

The system model includes two active RIS: one fixed on a building equidistant from all grids $\text{R}_G$, and the other is located on a mobile UAV $\text{R}_U$ which changes its location to provide maximum sum rate to users in the network. Both RISs serve all users in the network while being operated by a microcontroller to alter the phase shifts. Active RIS $R$, where $\text{R} = \{\text{R}_G, \text{R}_U\}$, in the model, is equipped with $k$ reflecting elements where $k \in \{1, 2, ..., K\}$. The signals reflected multiple times by the RISs are assumed to exhibit minimal power due to significant path loss as considered by the multi-RIS environment in [9], and this work carries the same assumption. Specifically, $\forall m \in \mathcal{M}, u \in \mathcal{U}$, and $o \in \mathcal{O}$, the positions of $\text{BS}_m$, $U_u$, and $O_o$ are represented by $p_m = (x_m, y_m, H_B)$, $p_u = (x_u, y_u, 0)$, and $p_o = (x_o, y_o, H_o)$, respectively, where $H_B$ and $H_O$ are the heights of the BSs and obstacles, respectively. The position of aerial RIS at time slot $t$ is denoted as $p_{R_U}[t] = (x_{R_U}[t], y_{R_U}[t], H_{R_U})$ and it hovers over a specific area $A$. As mentioned earlier, the UAV is mobile and changes its horizontal position in the xy-plane, while maintaining a fixed altitude $H_{R_U}$.

### B. Channel Model

The model includes line-of-sight and non-line-sight paths, so Rayleigh and Rician fading models are used. Due to scattering in the environment, Rayleigh fading channel models are used to simulate channels between the BSs and users. The channel
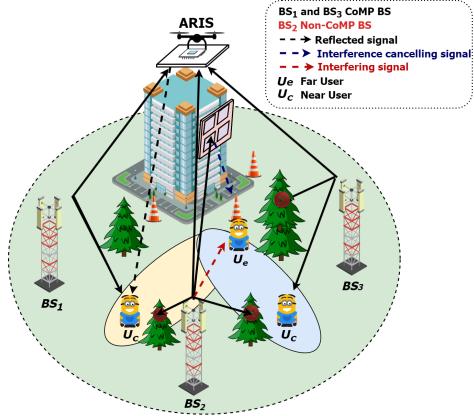


Fig. 1: RIS-assisted TN-NTN with coordinated NOMA pairing.

between $\text{BS}_m$ and user $\text{U}_u$ is denoted as $h_{m,u}$ and expressed for a time interval $t$ as

$$h_{m,u}[t] = \sqrt{\frac{\rho_0}{PL(d_{m,u})}} v_{m,u}[t], \qquad (1)$$

where $\rho_0$ is a reference path loss at 1 m, $PL(d_{m,u})$ is the large-scale path loss modeled as $PL(d_{m,u}) = (d_{m,u})^\alpha$, such that $\alpha$ is the path loss exponent, $d_{m,u}$ is the distance between $\text{BS}_m$ and $U_u$, whereas the small-scale Rayleigh fading coefficient is $v_{m,u} \sim \mathcal{R}(1)$. Since a direct line of sight exists between the BS and the RIS on the UAV and RIS on the building, this channel is simulated using Rician fading. The channel between the $\text{BS}_m$ and $R$ is denoted as $h_{m,R}$ and expressed for time interval $t$ as

$$h_{m,R}[t] = \sqrt{\frac{\rho_0}{PL(d_{m,R})}} \left( \sqrt{\frac{\kappa}{1+\kappa}} g_{m,R}^{\text{LoS}} + \sqrt{\frac{1}{1+\kappa}} g_{m,R}^{\text{NLoS}} \right), \tag{2}$$

where $\kappa$ is the Rician factor, $g_{m,R}^{\text{LoS}} \in \mathbb{C}^{K \times 1}$ is the LoS channel vector given by

$$g_m^{\text{LoS}} = \left[ 1, \ldots, e^{j(k-1)\pi \sin(\omega)}, \ldots, e^{j(K-1)\pi \sin(\omega)} \right]^T,$$

where $\omega$ represents the angle of arrival of the LoS component at $R$ while $g_{m,R}^{\text{NLoS}} \in \mathbb{C}^{K \times 1}$ is the NLoS component which follows Rayleigh fading as previously described [3]. This work assumes perfect channel state information (CSI) to avoid unnecessary overheads.

### C. Active RIS Configurations

The active RIS in the network amplifies the incoming signal, resulting from the reflection-type amplifiers employed by the active RIS elements backed by a power supply. The outgoing signal from the active RIS with $K$ reflecting elements is expressed as:

$$y[t] = \underbrace{P[t]\Theta[t]x[t]}_{\text{Desired signal}} + \underbrace{P[t]\Theta[t]v}_{\text{Dynamic noise}} + \underbrace{n_s}_{\text{Static noise}}, \qquad (3)$$

where $P[t] = \text{diag}(p_1[t], \ldots, p_K[t]) \in R^{K \times K}$ represents the amplification matrix of the active RIS for a time slot $t$. Each active RIS element can be denoted as $1 \leq p_N \leq \mathcal{S}$

where $\mathcal{S}$ is the maximum amplification an element can provide. Further, $x[t] \in \mathbb{C}^K$ denotes the incoming RIS signal, $y[t] \in \mathbb{C}^K$ denotes the outgoing RIS signal and $\Theta[t] = \text{diag}\left(a_1 e^{j\theta_1[t]}, a_2 e^{j\theta_2[t]}, \ldots, a_k e^{j\theta_K[t]}\right) \in \mathbb{C}^{K \times K}$ for a time instance $t$, where $a_k \in (0, 1]$ denotes the amplitude coefficient and $\theta_k \in [-\pi, \pi])$ denotes the phase shift of the $k$-th RIS element. Also, $diag(.)$ is the diagonalization operation that is performed on the RIS phase shift matrix. Active RIS elements consume additional power to amplify the reflected signals while generating thermal noise divided into two components: dynamic and static. Here, we assume only the dynamic noise and ignore static noise, since the study in [7] shows that static noise is disproportionately lesser in magnitude than dynamic noise and can be neglected. The dynamic noise variable is denoted as $v$ and modeled as $v \sim \mathcal{CN}(0_K, \sigma^2 I_K)$, where $\mathcal{CN}(\mu, \Sigma)$ signify the complex multivariate Gaussian distribution with mean and variance as $\mu$, $\Sigma$, respectively. The parameters $I_K$ and $0_K$ denote a $K \times K$ identity matrix and a $K \times 1$ zero vector, respectively.

### D. NOMA Model

All BSs serve a centre user and an edge user, so by following the 2-user NOMA tradition the power of the BS is distributed among the two users where the edge user receives a higher power signal [10], [11]. The signal transmitted by the $BS_m$ can be expressed as $x_m[t] = \sqrt{(1 - \lambda_m)p_t} x_{m,c}[t] + \sqrt{\lambda_m p_t} x_{m,e}[t]$, where $x_{m,c}[t]$ and $x_{m,e}[t]$ is the received signal for $U_c$ and $U_e$. All BSs in the system have the same transmit power, $p_t$, and the power allocation factor of $U_e$ is denoted as $\lambda_m$. The power allocation factor $\lambda$ should satisfy the constraint $0.5 < \lambda < 1$, since the edge user is further from the BS and demands higher power allocation [11].

Each $U_c$ receives a direct signal from the BS, an indirect reflected link from the RIS, interference from another BS serving a user of another cell, and inherent noise of active RIS. Ultimately, the signal received by $U_c$ is expressed as

$$y_c[t] = \underbrace{H_{m,c}[t]p_t x_m[t]}_{\text{desired signal}} + \underbrace{\sum_{j \neq m}^{M} H_{j,c}[t]p_t x_j[t]}_{\text{inter-user interference}} \qquad (4)$$
$$+ \underbrace{h_{R,c}[t]P[t]\Theta[t]v}_{\text{active RIS noise}} + n_u,$$

where $n_u$ denotes the additive white Gaussian noise for user $u$ as $n_u \sim \mathcal{N}(0, \sigma^2)$ and the channel from $BS_m$ to the user $U_c$ is defined as $H_{m,c}[t] = h_{m,c}[t] + h_{R,c}[t]P[t]\Theta[t]h_{m,R}$, which includes the direct link and the indirect link involving RIS. The inter-user interference is ICI obtained when $U_c$ experiences a signal broadcast from $BS_j$ [12]. After utilizing SIC, $U_c$ decodes the signal for $U_e$ and removes it from the signal received to decode its own signal. Therefore, the achievable rate for $U_c$ to decode the signal for $U_e$ is as follows

$$\mathcal{R}_{c \to e}[t] = \log\left(1 + \frac{\lambda_m \gamma_m}{(1 - \lambda_m)\gamma_m + 1}\right), \qquad (5)$$

where $\gamma_m = \frac{p_t |h_{m,c}[t]|^2}{\Psi}$ and $\Psi = p_t |h_{j,c}[t]|^2 + \sigma^2$. Whereas, the data rate for $U_c$ to decode its own signal is as follows

$$\mathcal{R}_c[t] = \log\left(1 + (1 - \lambda_m)\gamma_m\right). \qquad (6)$$

The signal received by $U_e$ is received via $BS_m$ and RIS $R$ can be expressed as

$$y_e = (h_{m,e}[t] + h_{R,e}^H[t]P[t]\Theta[t]h_{m,R}[t])x_m[t]p_t + (h_{j,e}[t] + h_{R,e}^H[t]P[t]\Theta[t]h_{j,R}[t])x_j[t]p_t + n_u, \qquad (7)$$

where $h_{R,e}^H$ represents the Hermitian transpose of the channel between RIS $R$ and user $e$. The data rate for $U_e$ can be expressed as

$$\mathcal{R}_e[t] = \log\left(1 + \frac{\lambda_m \gamma_m + \lambda_j \gamma_j}{(1 - \lambda_m)\gamma_j + (1 - \lambda_j)\gamma_j + 1}\right) \qquad (8)$$

After obtaining sum rate expression for centre user in (6) and for edge user in (8), we can obtain sum rate expression for whole system as

$$R_{\text{total}}[t] = \sum_{c \in \mathcal{C}} R_c[t] + \sum_{e \in \mathcal{E}} R_e[t]. \qquad (9)$$

### E. Energy efficiency

The energy efficiency of our proposed active RIS-assisted NOMA-CoMP system highlights the trade-off between network performance and energy costs, as follows:

$$\eta_E = \frac{R_{\text{total}}[t]}{p_{\text{total}}}. \qquad (10)$$

The total power consumed in the network comprises power from BS, active RIS, and UAV. The total sum rate of the system was expressed in (9), we can express the total power consumed in the network as follows

$$p_{\text{total}} = \underbrace{p_{\text{BS}}}_{\text{BS power}} + \underbrace{\sum_{i=1}^{2}\sum_{k=1}^{K}\left(\eta_{R_i}|P_{R_ik}\Theta_{R_ik}x_{R_ik}|^2 + P_{c,R_ik}\right)}_{\text{RIS power}} + \underbrace{P_{\text{UAV}}}_{\text{UAV power}}, \qquad (11)$$

where $p_{BS}$ represents the static power consumption and transmission power at the base station. The RIS power consumption includes $\eta_{R_i}$ as the amplification efficiency for the RIS $R_i$, $P_{R_ik}$ as the incoming signal power for the $k$-th element of the RIS $R_i$, $\Theta_{R_ik}$ as the phase shift applied by the $k$-th element of RIS $R_i$ panel, $x_{R_ik}$ as the incoming signal for the $k$-th element, and $P_{c,R_ik}$ denotes the circuit power consumption for the $k$-th element of RIS $R_i$. Lastly, $P_{UAV}$ indicates the power consumed by the UAV for its hovering and motion.

### F. Problem Formulation

In this work, our primary objective is to maximize the sum rate achieved over $t$ time slots. To achieve this goal, we jointly optimize three key control variables: the UAV trajectory denoted as $\widehat{\mathbf{T}} \triangleq \{p_{R_U}[t], \forall t\}$, the RIS phase shifts represented by $\mathbf{\Theta} \triangleq \{\Theta[t], \forall t\}$, the power allocation factors denoted as $\mathbf{\Lambda} \triangleq \{\lambda_m, \forall m\}$, and the RIS amplification matrix $\mathbf{P} \triangleq \{p_k[t], \forall k\}$ The problem can be mathematically formulated as

$$\max_{\mathbf{P},\boldsymbol{\Theta},\boldsymbol{\Lambda},\widehat{\mathbf{T}}} \sum_{t\in\mathcal{T}} R_{\text{sum}}[t], \qquad (10\text{a})$$

subject to:

$$x_R[t], y_R[t] \in A, \quad \forall t \in \mathcal{T}, \qquad (10\text{b})$$

$$\|\mathbf{p}_R[t] - \mathbf{p}_o\| \geq d_{\min}, \quad \forall o \in \mathcal{O}, t \in \mathcal{T}, \qquad (10\text{c})$$

$$\theta_k[t] \in [-\pi, \pi), \quad \forall k \in K, t \in \mathcal{T}, \qquad (10\text{d})$$

$$R_c[t] \geq R_c^{\min}, \quad \forall m \in M, t \in \mathcal{T}, \qquad (10\text{e})$$

$$R_e[t] \geq R_e^{\min}, \quad \forall e \in \mathcal{E}, t \in \mathcal{T}, \qquad (10\text{f})$$

$$\lambda_m \in (0.5, 1), \quad \forall m \in M, t \in \mathcal{T}, \qquad (10\text{g})$$

$$p_k[t] \in [1, \mathcal{S}], \quad \forall k \in K, t \in \mathcal{T}, \qquad (10\text{h})$$

## III. DEEP REINFORCEMENT LEARNING-BASED PROPOSED SOLUTION

### A. MDP Formulation

The MDP is defined via the tuple $(\mathbb{S}, \mathbb{A}, \mathbb{P}, \nabla, \mathbb{D})$, where $\mathbb{S}$ and $\mathbb{A}$ represent the state space and action space, respectively. $\nabla$ is the reward to the agent for its actions, and $\mathbb{D}$ is the discount factor to balance weights of immediate and future rewards. Lastly, $\mathbb{P}$ denotes the state transition probability i.e., the likelihood of transitioning from one state to another for an action. At each time slot $t$, the agent observes the current state $s_t$, selects an action $a_t$ based on its policy, and transitions to a new state $s_{t+1}$, to get rewarded.

### B. State Space

The state space in a time slot $t$ is composed of the current position of UAV, BS power allocation factor, amplification matrix of active RIS, and achievable sum rates, which can be denoted as $\varrho[t]$, $\zeta$, $\vartheta[t]$, and $\mathbf{R}[t] = \{\mathcal{R}_c[t], \mathcal{R}_e[t], \forall c, e\}$, respectively. Eventually, the state space becomes

$$s_t = \{\varrho[t], \zeta, \vartheta[t], \mathbf{R}[t]\}. \qquad (12)$$

The dimension for the state space can be expressed as $\dim(\mathbb{S}_t) = 2 + M + K^2 + 2$.

### C. Action Space

The action space for the MDP is composed of the UAV's horizontal movement in xy-plane, the amplification factor and phase shifts of the active RIS, and the power allocation factor of the BSs. Specifically, the action space at time slot $t$ contains the UAV actions $a_R[t] \in \{(-1,0), (1,0), (0,-1), (0,1), (0,0)\}$, representing left, right, down, up, and hover action, respectively. The phase shifts $a_\theta[t] = \{\theta_k[t], \forall k\}$, the power allocation factors $a_\Lambda = \{\lambda_m, \forall m\}$, and the active RIS amplification matrix $a_p[t] = \{p_k, \forall k\}$. Thus, the action space can be expressed as

$$a_t = \{a_R[t], a_\theta[t], a_\Lambda, a_p[t]\} \qquad (13)$$

The dimension for the action space is expressed as $\dim(a_t) = 2 + K + M + K^2$.

---

**Algorithm 1** H-PPO Active RIS-Based Energy Optimization
1: **Initialize:** Parameters $\alpha_d, \alpha_c, \zeta, \Psi$
2: **for** iteration $l = 1, \ldots, L$ **do**
3:     Receive initial state $s_0$
4:     **for** time step $t = 0, \ldots, T$ **do**
5:         Choose continuous actions $a_p$ and $a_c$ using $\pi_{\alpha_c}(s_t)$
6:         **for** each RIS element $n$ **do**
7:             Energy optimization $e_n = g(\psi_n, a_p, a_c)$
8:             Update the gain for the $k$-th RIS element based on channel conditions
9:         **end for**
10:        Execute actions $a_t = \{a_R, \Psi, a_p, a_c\}$
11:        Observe reward $r_t$ based on achieved sum rate and outage probability
12:        Store the transition $(s_t, a_t, r_t, s_{t+1})$ in experience replay buffer
13:        Compute the advantage function estimate $\hat{A}_t$
14:    **end for**
15: **end for**
16: **for** epoch $m = 1, \ldots, K$ **do**
17:     Sample mini-batch experiences $E$ from the replay buffer
18:     Objective computation using clipped functions for both discrete and continuous actions: $L_d^{\text{CLIP}}(\theta_d)$ and $L_c^{\text{CLIP}}(\theta_c)$
19:     Optimize policy parameters $\alpha_d$ and $\alpha_c$
20: **end for**
21: Obtain previous policy parameters $\alpha_d^{\text{old}} \leftarrow \alpha_d$ and $\alpha_c^{\text{old}} \leftarrow \alpha_c$
22: Clear

---

### D. Reward Function

The reward function ensures the maximum sum rate while ensuring UAV safety and meeting QoS requirements by penalizing constraint violations when UAV goes OOB (out of bounds). The reward function is defined as

$$R(s_t, a_t) = R_{\text{sum}}[t] + \xi_{\text{dist}} \left( \frac{C}{d_{\text{R}_U,\text{U}}[t]} \right) \zeta - \xi_{\text{OOB}} \cdot \mathbb{I}(\text{OOB}), \quad (14)$$

where $\zeta[t] = \mathbb{I}(d_{\text{R}_U,\text{U}}[t] < \text{threshold})$ acts as the indicator function for the distance incentive to keep the UAV close to users. Also, $\xi_{\text{OOB}}$ represents the penalty given to the agent when $\text{R}_U$ goes out of bounds of the grid which is pointed by the indicator function $\mathbb{I}(\text{OOB})$. Lastly, $\xi_{\text{dist}}$ and $C$ are kept as constants in the expression.

### E. H-PPO algorithm

We employ a hybrid Proximal Policy Optimization (H-PPO) algorithm to cater to the hybrid action space in our network. The algorithm offers actions in both discrete $a_R[t]$ and continuous action spaces $a_\theta[t], a_\Lambda, a_p[t]$. The value function $V(\mathbb{S}_t)$ is used to obtain a variance-reduced advantage function estimate $\hat{A}_t$ to optimize policy. Following the implementation details used in [13], the policy is executed for $T$ time steps, and $\hat{A}_t$ is computed as

$$\hat{A}_t = \sum_{k=0}^{\bar{T}-1} \mathbb{D}^k r_{t+k} + \mathbb{D}^{\bar{T}} \hat{V}(s_{t+\bar{T}}) - V(s_t), \qquad (15)$$

where $\bar{T}$ is much smaller than the length of the episode $T$. Stochastic policy for discrete and continuous action is

generated in the same way so only discrete is mentioned here. To generate the stochastic policy $\pi_{\theta_d}(a_t|s_t)$ for the discrete actions, the corresponding actor network outputs $|\mathbf{a}_R|$ logits, which are then passed through a softmax function to obtain a probability distribution over the available discrete actions. Conversely, the continuous actor network generates the continuous actions $a_\Phi$ and $a_\Lambda$ by sampling from Gaussian distributions parameterized by the mean and standard deviation outputs of the network, as dictated by the stochastic policy $\pi_{\theta_c}(a_t|s_t)$. Both $\pi_{\theta_d}(a_t|s_t)$ and $\pi_{\theta_c}(a_t|s_t)$ are optimized independently using their respective clipped surrogate objective functions. For the discrete actions, the objective function is given by

$$L_d^{\text{CLIP}}(\theta_d) = \hat{\mathbb{E}}_t \left[ \min \left( r_t^d(\theta_d)\hat{A}_t, \aleph(r_t^d, \theta_d, \epsilon)\hat{A}_t \right) \right] \quad (16)$$

where $\aleph(r_t^d, \theta_d, \epsilon) = \text{clip}(r_t^d(\theta_d), 1 - \epsilon, 1 + \epsilon)$, $r_t^d(\theta_d) = \frac{\pi_{\theta_d}(a_t|s_t)}{\pi_{\theta_d^{\text{old}}}(a_t|s_t)}$ is the importance sampling ratio, and $\epsilon$ is the clipping parameter.

The optimization objectives of both policies remain decoupled i.e., $\pi_{\theta_d}(a_t|s_t)$ and $\pi_{\theta_c}(a_t|s_t)$ are treated as independent distributions during policy optimization, rather than a joint distribution encompassing both action spaces. The H-PPO algorithm is summarized in Algorithm 1.

## IV. SIMULATION RESULTS

### A. Simulation setup

In order to assess the effectiveness of the proposed active RIS framework in a CoMP-NOMA network, we create a virtual representation of an urban setting which includes $C = 3$ BSs, $K = 3$ users, and $L = 0$ obstacles. The initial position of the UAV is set at $(-5, 0, 40)$ m, while the BSs are positioned at coordinates $[-30, 30, 20]$, $[30, 30, 20]$, and $[20, -30, 20]$, respectively. Each base station, however, is considered to operate at a unique power level, denoted as $P_{d,x}$ for each base station $x$, where $x \in \{1, 2, 3\}$. The remaining utilities, users inclusive, are distributed randomly throughout the environment. The network operates at a carrier frequency of $f_c = 2.4$ GHz, utilizing a bandwidth of $B = 10$ MHz. The noise power is defined as $\sigma^2 = -174 + 10\log_{10}(B)$, dBm. To account for attenuation effects, path loss exponents are defined as $\beta_n = 2.2$, $\beta_o = 3.0$, $\beta_k = 3.3$, and $\beta_s = 3.7$. The results presented are averaged over 1000 independent realizations of user positions to ensure reliable performance metrics.

### B. Numerical results

In Fig.2a investigation of the relationship between transmit power $p_t$, where $p_t \in [-30\,\text{dBm}, 30\,\text{dBm}]$ and sum rate $R_k$ is made, comparative analysis between different H-PPO-RIS configuration were taken in to account, highlighting active RIS improvement. Compared with H-PPO P-RIS $R_k$, H-PPO A-RIS achieves $(7\,\text{bps/Hz})$ at $p_t$ of $-10$ dBm, performing better than H-PPO P-RIS configuration. The cumulative rate $R_k(T, P_{\text{tx}}, h, \gamma)$ for H-PPO P-RIS and H-PPO A-RIS increases to $(15\,\text{bps/Hz})$ and $(20\,\text{bps/Hz})$, respectively, when $p_t$ reaches 25 dBm. This represents a 33% improvement in sum

rate, emphasizing the amplification capability of active RIS enabled by NOMA, which outperforms passive configurations..

Our investigation further spread to the learning algorithm efficacy for energy optimization. Fig.2b illustrates the progression of reward dynamics $r$ across multiple training steps, with HPPO-NOMA setup displaying a steeper growth, from 0.2 to 1.0, showing higher learning efficiency, indicating NOMA's improved energy management.

Also, we have evaluated the outage probabilities of our configurations, In Fig.2c at $p_t$ of $-30$ dBm to $-10$ dBm, outage probability of all configuration remains same for all worst users rate. The HPPO A-RIS NOMA configuration has the least outage probability 0.20 at 30 dBm, indicating the worst user rate could increase from 5 bps/Hz, hence outperforms all other configurations in consideration.

Fig.3a illustrates the relationship between energy efficiency and spectral efficiency for various configurations. The H-PPO A-RIS NOMA configuration demonstrates the highest performance, reaching a peak energy efficiency of $11,500$ bit/J at 24 bits/s/Hz.This improvement results from the active RIS's ability to dynamically adjust reflections, enhanced further by the NOMA strategy, achieving a 20% improvement over passive RIS configurations. Fig.3b shows the sum rate for different configurations as the number of RIS elements increases. $K$ increases from 30 to 200 for all configurations indicating that adding more RIS elements $K$ enhances the system's ability to optimize reflections and HPPO A-RIS NOMA consistently shows the highest $R_k$ across all $K$ highlighting the advantage of active RIS with NOMA.

Lastly, we have Fig.3c illustrating the relationship between $p_t$ per base station and $R_k$ for different RIS configurations with fairness and without fairness, indicating the trade-off between system capacity and fairness among users. As the $p_t$ increases from $-5$ dBm to 25 dBm, the $R_k$ improves for all configurations. With an even power distribution among center and edge users, fairness slightly reduces the $R_k$. H-PPO A-RIS-NOMA with fairness achieves the highest performance among all other configurations.

## V. CONCLUSION

This paper proposes a DRL-based resource allocation framework for active RIS-assisted TN-NTN networks using CoMP-NOMA.The active RIS response system incorporating COMP-NOMA significantly achieved higher sum rates and energy efficiency than the traditional passive RIS and other OMA configurations. These results show that the active RIS integration with CoMP-NOMA and DRL is an efficient option for resource allocation and better system performance. Future work could explore the integration of energy harvesting technologies with active RIS, or using mmWave and terahertz communication to open exciting research avenues.

## VI. ACKNOWLEDGEMENTS

(a) Sum Rate vs transmit power.

(b) Reward function vs time steps with normalized RL policies
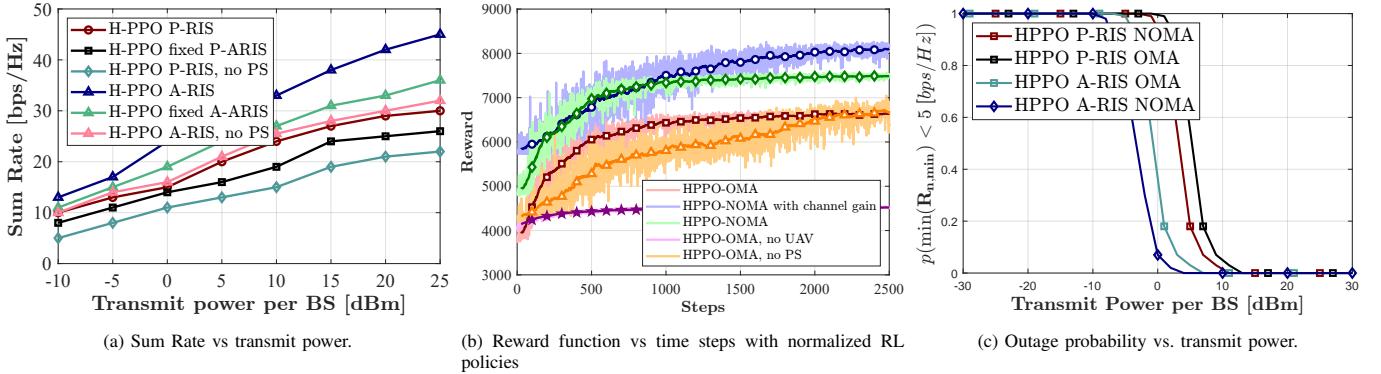
(c) Outage probability vs. transmit power.

Fig. 2: **(a)** Sum rate performance for baseline methods [OMA, Brute Force and RL-NOMA] varying transmit power, **(b)** Training convergence of reward functions for various RL configurations over timesteps, **(c)** Outage probability comparison against transmit power per base station for NOMA and OMA base.
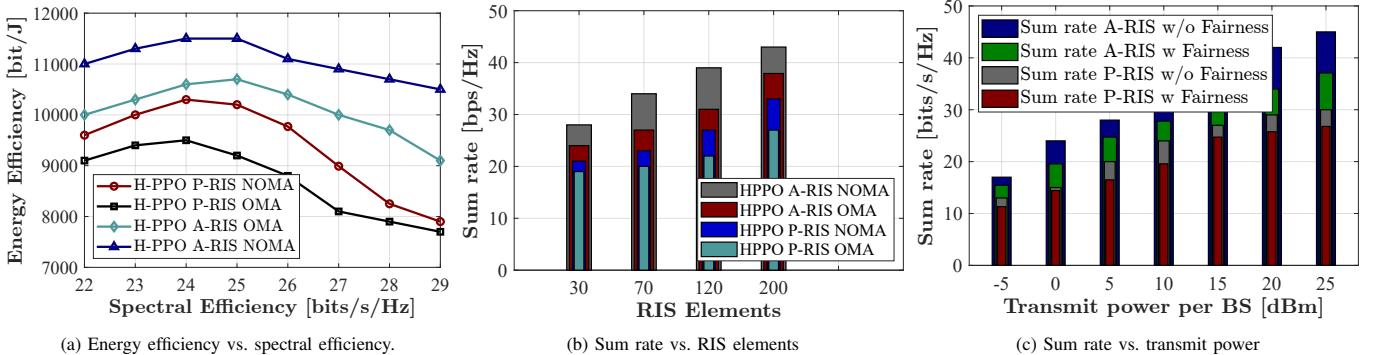


(a) Energy efficiency vs. spectral efficiency.

(b) Sum rate vs. RIS elements

(c) Sum rate vs. transmit power

Fig. 3: **(a)** Energy efficiency comparison against spectral efficiency between different RIS configuration , **(b)** Sum rate comparison with RIS elements n belong to [30,70,120,200] for different RIS configuration **(c)** Sum rate comparison against baseline methods [DRL NOMA (w and w/o fairness), OMA].

## REFERENCES

[1] M. Harounabadi and T. Heyn, "Toward integration of 6G-NTN to terrestrial mobile networks: Research and standardization aspects," *IEEE Wireless Communications*, vol. 30, no. 6, pp. 20–26, 2023.

[2] Y. Liu, X. Mu, X. Liu, M. Di Renzo, Z. Ding, and R. Schober, "Reconfigurable intelligent surface-aided multi-user networks: Interplay between NOMA and RIS," *IEEE Wireless Communications*, vol. 29, no. 2, pp. 169–176, 2022.

[3] M. Umer, M. A. Mohsin, S. A. Hassan, H. Jung, and H. Pervaiz, "Performance analysis of STAR-RIS enhanced CoMP-NOMA multi-cell networks," in *2023 IEEE Globecom Workshops (GC Wkshps)*, 2023, pp. 2000–2005.

[4] Y. Ge and J. Fan, "Active reconfigurable intelligent surface assisted secure and robust cooperative beamforming for cognitive satellite-terrestrial networks," *IEEE Transactions on Vehicular Technology*, vol. 72, no. 3, pp. 4108–4113, 2022.

[5] R. H. Ratul, M. Iqbal, T. Ashraf, J.-Y. Pan, Y.-H. Wang, and S.-Y. Lien, "Adaptive three layer hybrid reconfigurable intelligent surface for 6g wireless communication: Trade-offs and performance," in *2023 IEEE Asia Pacific Conference on Wireless and Mobile (APWiMob)*, 2023, pp. 232–236.

[6] H. Mei, K. Yang, Q. Liu, and K. Wang, "3D-Trajectory and Phase-Shift Design for RIS-Assisted UAV Systems Using Deep Reinforcement Learning," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 3, pp. 3020–3029, 2022.

[7] Z. Zhang, L. Dai, X. Chen, C. Liu, F. Yang, R. Schober, and H. V. Poor, "Active RISs: Signal modeling, asymptotic analysis, and beamforming design," in *GLOBECOM 2022 - 2022 IEEE Global Communications Conference*, 2022, pp. 1618–1624.

[8] S. Ammar, C. P. Lau, and B. Shihada, "An in-depth survey on virtualization technologies in 6g integrated terrestrial and non-terrestrial networks," *IEEE Open Journal of the Communications Society*, 2024.

[9] M. Bilal, S. F. Zahra, H. Rizwan, T. Umar, S. A. Hassan, H. Jung, and K. Dev, "Single versus double IRS-assisted networks: A comparative analysis using practical phase shifting," in *2024 IEEE Wireless Communications and Networking Conference (WCNC)*, 2024, pp. 1–5.

[10] M. Elhattab, M. A. Arfaoui, C. Assi, and A. Ghrayeb, "RIS-assisted joint transmission in a two-cell downlink NOMA cellular system," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 4, pp. 1270–1286, 2022.

[11] M. Obeed, H. Dahrouj, A. M. Salhab, S. A. Zummo, and M.-S. Alouini, "User pairing, link selection, and power allocation for cooperative NOMA hybrid VLC/RF systems," *IEEE Transactions on Wireless Communications*, vol. 20, no. 3, pp. 1785–1800, 2021.

[12] Z. Shi, H. Lu, X. Xie, H. Yang, C. Huang, J. Cai, and Z. Ding, "Active RIS-aided EH-NOMA networks: A deep reinforcement learning approach," *IEEE Transactions on Communications*, vol. 71, no. 10, pp. 5846–5861, 2023.

[13] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *Proceedings of The 33rd International Conference on Machine Learning*, M. F. Balcan and K. Q. Weinberger, Eds., vol. 48, 2016, pp. 1928–1937.