

# Multi-agent Attention Actor-Critic Algorithm for Load Balancing in Cellular Networks

Jikun Kang, Di Wu, Ju Wang, Ekram Hossain, Xue Liu, Gregory Dedek  
Samsung Electronics, Canada

{jikun.kang, e.hossain}@partner.samsung.com, {di.wu1, j.wang1, steve.liu, greg.dudek}@samsung.com

**Abstract**—In cellular networks, User Equipment (UE) handoff from one Base Station (BS) to another, giving rise to the load balancing problem among the BSs. To address this problem, BSs can work collaboratively to deliver a smooth migration (or handoff) and satisfy the UEs’ service requirements. This paper formulates the load balancing problem as a Markov game and proposes a Robust Multi-agent Attention Actor-Critic (Robust-MA3C) algorithm that can facilitate collaboration among the BSs (i.e., agents). In particular, to solve the Markov game and find a Nash equilibrium policy, we embrace the idea of adopting a nature agent to model the system uncertainty. Moreover, we utilize the self-attention mechanism, which encourages high-performance BSs to assist low-performance BSs. In addition, we consider two types of schemes, which can facilitate load balancing for both active UEs and idle UEs. We carry out extensive evaluations by simulations, and simulation results illustrate that, compared to the state-of-the-art MARL methods, Robust-MA3C scheme can improve the overall performance by up to 45%.

**Index Terms**—Load balancing, multi-agent reinforcement learning, Nash equilibrium, self-attention

## I. INTRODUCTION

In a mobile cellular system, User Equipment (UE) movements across Base Stations (BSs) result in a Load Balancing (LB) problem causing UE dissatisfaction [1]. To balance the load, UEs are migrated among BSs and channels in each BS. Thus, it is crucial to coordinate multiple base stations to reallocate UEs when they move from one region to another. In a 5G network, BSs could be more close to each other compared with LTE/4G network, which means UEs are more likely to move from one BS to another more frequently. This makes BSs’ collaboration even more important than before.

Recently, Reinforcement learning (RL) approaches have shown the advantages of solving the load balancing problem [2]. However, existing RL-based LB algorithms cannot address multi-BSs LB problem properly. The single agent solution is proposed to control BSs with one policy [3]. Hierarchical [4] and transfer learning-based RL methods [5] have also been developed and showed improved performance in terms of data throughput and load variation reduction. However, these methods fail to consider multi-agent interactions and cannot guarantee Nash equilibrium. In addition, the classical multi-agent actor-critic algorithm [6] does not consider Nash equilibrium.

To address aforementioned issues, in this paper, we formulate the load balancing problem over multiple base stations (Fig. 1) as a Markov game. Specifically, we treat base stations

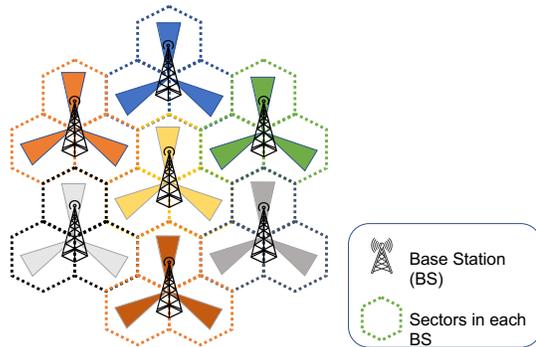


Fig. 1: A cellular network tessellation with 7 BSs. Each hexagon denotes one sector. One BS consist of three sectors, which draw in the same color. Each sector controls 120 angle degrees. Furthermore, in each sector, there are 4 channels residing on 4 different carrier frequencies. In this paper, we control the parameters of these 7 BSs to achieve load balancing.

as agents and learn the corresponding policies to control the load balancing. To solve the Markov game, we propose Robust-MA3C, a Robust-Multi-agent Attention Actor-Critic algorithm. The main focus of our work is to leverage collaboration among BSs to improve system performance. In particular, we seek to enhance the UE experience in under-performing BSs. A major challenge in achieving these goals is that agents tend to maximize their own goals greedily, which affect the overall experiences of the UEs. The state-of-the-art rule-based method cannot address this challenge properly, since they do not formulate the problem in multi-agent settings. The single-agent approach available in the literature is not suitable to foster BS collaboration [3]. This is because, in single-agent algorithms, BSs are controlled by one centralized neural network, which fails to converge when the action space increases with the number of agents. The state-of-the-art MARL load balancing algorithms are not able to deal with this challenge either, since they spare no effort to preserve the UE experience in poorly performing BSs.

To overcome this challenge, we adopt the following two mechanisms. First, to encourage collaboration among all BSs, we introduce an attention-based [7] multi-agent actor critic algorithm. Specifically, during the centralized training, each transmitted information will be assigned a weight value. The

motivation is that agents will have different attention to the received information. Since our objective is to maximize the overall system performance, high-performance agents can help the poorly performing agents based on transmitted information. Second, we utilize the nature agent idea [8] to achieve the Nash equilibrium policy. The nature agent plays against all other agents by selecting the worst-case actions. Thus, to fight against the worst-case scenario raised by the nature agent, all agents need to work together and develop a joint equilibrium policy.

We evaluate the proposed Robust-MA3C method against the state-of-the-art (SOTA) multi-agent RL (MARL)-based methods in a system-level network simulator for a network, which is depicted in Fig. 1. The simulation results show that our method can improve the network throughput by up to 45%.

In summary, this paper makes the following contributions:

- We formulate the multi-BS LB problem as a Markov game and propose a MARL solution, which enhances system performance by exploiting collaborations among the BSs.
- To improve the underperforming BSs' performance, we propose an attention-based weighting mechanism to each transmitted information.
- To model uncertainty, we add a nature agent to play against all agents to find the Nash equilibrium policy.
- We benchmark the proposed method against five other methods and demonstrate its performance superiority.

## II. RELATED WORK

Load balancing algorithms have been designed for wireless networks to evenly distribute UEs across different network resources, i.e., frequency bands (channels) and base stations. The MAS model [9] is for a dynamic load balancing scheme for a multi-agent system in which the agent selection is based on agents' credit value, and location selection is based on the inter-machine communication load. Another mobility-based load balancing infrastructure is called JIAC [10]. The migration decision is taken locally which prevents the centralized control problems. In [11], a centralized load balancing system is given. The selection policy is based on the job's execution time, while the location policy is based on negotiation with cluster nodes. Previous work [3] has explored single-agent RL for load balancing in network systems, with a close-to-heuristic performance achieved in moderate-scale simulations. This paper applies and studies MARL algorithms on the network load balancing problem in a realistic testbed, which considers load balancing for idle UEs and active UEs at the same time.

Approaches that address multi-agent cooperative tasks can be generally grouped into two types. The first type is based on implicit message exchange. These types of methods utilize a centralized training framework to share agents' value function parameters, which facilitates agents' collaboration. MADDPG [12] adopts global observations and actions for each action-value function estimation. COMA [13] leverages the advantage of this idea and estimates each advantage function using a

counterfactual baseline. The second type of methods is based on explicit communication using learned protocols [14]–[16]. In DIAL [14], each agent generates message and Q-value function at the same time. Each agent then encoded the learned message and computed it with other transmitted information. CommNet [16] learns a generalized communication protocol, which aggregates all of the received messages and computes the average value. NeurComm [15] propose a differentiable communication protocol, which can reduce information loss. In this paper, we focus on the first type with implicit message exchange, because it is more applicable for a real-world scenario without additional communication protocol requirements.

## III. SYSTEM MODEL AND ASSUMPTIONS

### A. Wireless Network Terminologies

We use the term “load” to refer to the number of UEs being served. The term “Base Station” (BS) describes a physical site, where radio access devices are placed. In each BS, there are three non-overlapping “sectors”, each of which controls 120 angle degrees. A sector serves the UEs located in a certain direction which is associated with the corresponding BS. Each sector contains several “channels” corresponding to the carrier frequencies supported by a sector. A channel is a service entity that serves UEs in a certain direction and on a specific carrier frequency.

### B. Communication Load Balancing Features

In this work, we focus on leveraging reinforcement learning to optimize load balancing for two types of users, i.e., idle UEs and active UEs. The idle UE load balancing (IULB) is a load balancing feature designed specifically for idle UEs. This load balancing technique can help adjust the host channels for idle UEs from a crowded channel to a less overloaded channel. The channel is referred to as a combination of the sector and frequency. The control knobs for IULB are the channel re-selection ratios, which lie between 0 and 1. By adjusting the re-selection ratios, the probabilities of frequency to be selected as the host channel will be adjusted. Then, when the status for UEs changes from idle to active, the communication load will be distributed more evenly between different cells. Active UE load balancing (AULB) is a load balancing feature designed for active UEs that are actively exchanging data with the BS. The AULB will be triggered when the load of the current serving channel is larger than the neighbouring channel, and one of the neighbouring channels will be chosen. For AULB, there will be two threshold values for each channel, which help to determine when the load balancing feature will be determined and which neighbouring channel will be determined as the target channel.

### C. System Metrics

Suppose that there are  $N_U$  UEs (either active, idle, or mobile) in the network. Define  $U_i$  as the set of UEs associated with the  $i$ -th channel. Among  $U_i$ , there are both active UEs (denoted as  $U_i^a$ ) and idle UEs (denoted as  $U_i^d$ ). Naturally, we

have  $U_i = U_i^a \cup U_i^d$ . Further, let  $u_{i,k}$  denote the  $k$ -th UE in the  $i$ -th channel. Note that an idle UE at the current moment may become active in the future, and vice versa. We aim to balance the assignments of UEs to different channels, to improve the system performance in terms of the following performance metrics.

**Average Throughput**  $G_{aver} = \frac{1}{N_U} \sum_i \sum_k \frac{A_{i,k}}{T}$  measures the overall system performance, where  $T$  is the time period of interest, and  $A_{i,k}$  denotes the total size of packets received by  $u_{i,k}$  within  $T$ . Improving this metric means increasing the overall system performance.

**Minimum Throughput**  $G_{min} = \min_{i,k} \left( \frac{A_{i,k}}{T} \right)$  captures the worst-case UE performance.

**Standard Deviation**  $G_{sd} = \left( \sqrt{\frac{1}{N_U} \sum_i \sum_k \left( \frac{A_{i,k}}{T} - G_{aver} \right)^2} \right)^{-1}$  represents the fairness services to all UEs. Minimizing this metric reduces the gap between different UEs' performance, and improves fairness to UEs.

We linearly combine these metrics as the reward  $r = G_{aver} + G_{min} - G_{sd}$ . Note that these metrics are in the same range, which does not require normalization.

#### IV. METHODOLOGY

##### A. Definitions and Problem Statement

There exist  $N$  BSs, which allocate UEs to different channels based on their current location and usage. All BSs work together to satisfy the UEs' communication requirements. In this work, we formulate this multi-base-station load balancing task as a Markov game  $G$ , which is comprised of a set of states and possible actions. Each action transits the current state to a new state and receives a reward. The objective is to find the optimal policy, i.e., an action in each state to maximize the expected total reward. A Markov game can be described by a tuple  $G = (N, S, A, P, R)$ :

- $N$ : the set of agents in the Markov game. Here, we denote them as  $N = \{1, \dots, N\}$ , where  $N$  is the total number of base stations.
- $S$ : the set of joint states. We denote agent  $k$ 's state as  $s^k$ . In our case, there are three types of states, which are number of UEs in every BS  $s_{ue}$ , the bandwidth utilization of every BS  $s_{band}$ , and the average throughput of every BS  $s_{tput}$ .
- $A$ : the set of joint actions of all agents,  $a = \{a^1 \times \dots \times a^N\}$ . Agent  $k$ 's action is denoted by  $a^k$ . In particular, every action  $a^k$  contains two components. The first component regulates the active UE load balancing operations and the second component controls the idle UE load balancing parameters.
- $P$ : the state transition probability function, where  $P(S_{t+1} = s' | S_t = s, A_t = a)$  maps a state-action pair at time  $t$  to a probability distribution over states at time  $t+1$ , such that  $P : S \times A \times S' \rightarrow [0, 1]$ .
- $R$ : the reward function. In the multi-agent cooperation setting, the goal is to maximize the total reward  $r =$

$r^1 + \dots + r^n$ , where  $R^k : S^k \times A^k \rightarrow R^k$  is the reward of agent  $i$  (to be defined in Sec. III-C).

Besides, we denote by  $\pi = \{\pi^1, \dots, \pi^N\}$  the set of policies of all agents. Formally, we define the multi-agent LB problem as follows:

$$\mathcal{J}(\pi) = \max_{\alpha^k, \beta^k, \gamma^k \sim \pi^k} \sum_{k=1}^N G_k, \quad (1)$$

$$\text{s.t. } \alpha^k \in [\alpha_{min}^k, \alpha_{max}^k], \quad (2)$$

$$\beta^k \in [\beta_{min}^k, \beta_{max}^k], \quad (3)$$

$$\gamma^k \in [\gamma_{min}^k, \gamma_{max}^k], \quad (4)$$

where  $G_k$  denotes the system performance of BS  $k$ ,  $\alpha_{min}^k$  and  $\alpha_{max}^k$  denote the controllable range of AULB actions of BS  $k$ , and  $\beta_{min}^k, \beta_{max}^k, \gamma_{min}^k$  and  $\gamma_{max}^k$  define the controllable range of IULB actions of BS  $k$ . In this paper,  $\alpha^k \in [-2db, 2db], \beta^k \in [-20db, 20db], \gamma^k \in [-20db, 20db], k = \{1, \dots, N\}$ . *Due to its combinatorial nature, solving this multi-agent LB problem using traditional optimization methods will not be feasible for a practical system.*

The goal of MARL is to learn a set of agent policies  $\{[\alpha_{i,j}^k, \beta_{i,j}^k, \gamma_{i,j}^k] \sim \pi^k\}_{k=1, \dots, N}$  that maximize the total expected return  $\sum_{k=1}^N G_k$  per episode.

For multi-agent load balancing, since we focus on maximizing the joint reward of all BSs, the proposed solution should be in accordance with the Nash equilibrium. In particular, the Nash equilibrium is defined as follows:

$$\mathcal{J}(\pi_*^k, \pi_*^{-k}) \geq \mathcal{J}(\pi_*^k, \pi_*^{-k}), \forall k \in N, \quad (5)$$

where  $\pi_*^{-k} = \prod_{k \neq i} \pi^i$  refers to other agents' joint policies except agent  $k$ . Our goal is to find a Nash equilibrium such that, given all other agents' equilibrium policies  $\pi_*^{-k}$ , there is no motivation for agent  $k$  to deviate from  $\pi_*^k$ .

##### B. Actor-Critic Algorithm

To start with, we first introduce the actor-critic algorithm, which is the primitive of the state-of-the-art algorithms in MARL [12], [13]. The actor-critic algorithm [17], as its name suggests, consists of two functions: the actor function and the critic function. The critic function estimates the action value, which is often noted as the Q-value function. This function serves as the policy learning guidance as it estimates the performance of the current action, which can be seen as an expert that controls the direction of gradient learning. The actor function directly controls the agent's behaviors according to the "suggestion" from a critic.

##### C. Multi-agent Actor-Critic Algorithm

We propose to extend the actor-critic algorithm to the multi-agent scenario. We adopt the centralized critic training decentralized actor execution framework [12], [13]. During the critic training phase, each agent receives neighbor agents' information (the size of the transmitted information is about a few bytes) at every step  $t$ , which are states  $s_i$  and actions  $a_i$ . Thus, the critic loss function can be written as:

$$\mathcal{L}(\phi^i) = \mathbb{E}_{s,a,r} [y - Q^i(s_t^1, a_t^1, \dots, s_t^k, a_t^k; \phi^i)]^2, \quad (6)$$

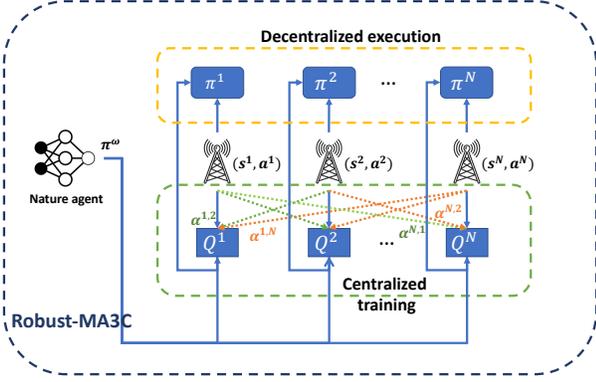


Fig. 2: Illustration of proposed algorithms. Dash lines refer to the information transmission. Different dash line colors correspond to different importance weights.

where  $y = r^i + \gamma Q^i(s_{t+1}^1, \mathbf{a}_{t+1}^1, \dots, s_{t+1}^k, \mathbf{a}_{t+1}^k; \phi^i)$ . While in the execution phase, the agent makes decisions based on its states. According to the policy gradient theorem [18], the actor updating function can be written as:

$$\nabla \mathcal{J}(\theta^i) = \mathbb{E}_{s^i, \mathbf{a}^i \sim \mathcal{D}} [\nabla_{\theta^i} \log \pi^i(\mathbf{a}^i | s^i) Q^i(s_t^1, \mathbf{a}_t^1, \dots, s_t^k, \mathbf{a}_t^k)],$$

#### D. Attention-Based Multi-agent Actor-Critic Algorithm

To encourage collaboration among different agents, we propose an attention-based multi-agent learning approach. Specifically, in the centralized critic training phase, we adopt the self-attention mechanism [7] to learn the importance weights of transmitted information. For the simplicity of calculation, we encode information  $e^i = g(s^i, \mathbf{a}^i)$ , where  $g(\cdot)$  is a multi-layer perceptron (MLP). Hence, the importance weight of agent  $j$  to agent  $i$  is denoted as  $\alpha^{i,j} = \text{softmax}\left(\frac{e^i \cdot e^j}{\sqrt{d_k}}\right) e^i$ , where  $d_k$  is the dimension of  $e^k$ . Then, we assign an attention weight to each received information  $e^j$  and re-write the critic loss function as follows:

$$\mathcal{L}(\phi^i) = \frac{1}{N} \sum_j (y - Q^i(s^i, \alpha^{i,1} e^1, \dots, \alpha^{i,N} e^N; \phi^i))^2, \quad (7)$$

where  $y = r^i + \gamma Q^i(e^i, \alpha^{i,1} e^1, \dots, \alpha^{i,N} e^N; \phi^i)$ ,

where  $Q^i(\cdot)$  is the centralized action-value function [12], to which we extend the neighbor attention. The inputs of this function are encoded by agents' observations and actions. The outputs are action-value for agent  $i$ .

Moreover, with the introduced attention weights, the actor updating function can be expressed as:

$$\nabla \mathcal{L}(\theta^i) = \mathbb{E}_{s^i, \mathbf{a}^i \sim \mathcal{D}} [\nabla_{\theta^i} \log \pi^i(\mathbf{a}^i | s^i) Q^i(e^i, \alpha^{i,1} e^1, \dots, \alpha^{i,N} e^N)]. \quad (8)$$

#### E. Improving Nash Equilibrium by Adding a Nature agent

To further enhance agents' collaborations and solve the Nash equilibrium defined in Eqn. 5, we adopt the nature agent idea [8] and propose a corresponding Robust-MA3C algorithm. The motivation is that the nature agent always plays against each agent by selecting the worst-case actions at every

state, which can be also seen as the model of the reward uncertainty. Thus, to fight against the worst-case scenario raised by the nature agent, all agents need to work together and develop a joint equilibrium policy.

We replace the reward with the nature agent in the critic function and develop a robust-critic loss function:

$$y = \pi^\omega(\hat{s}, \hat{\mathbf{a}}) + \gamma Q^i(e^i, \alpha^{i,1} e^1, \dots, \alpha^{i,N} e^N; \phi^i), \quad (9)$$

where nature-agent's policy is approximated by an neural network parameterized by  $\omega$ ,  $\hat{s}$  stands for the nature-agent's state, and  $\hat{\mathbf{a}}$  refers to nature-agent's action. We update the nature agents policy parameter as follows:

$$\omega_{t+1} = \omega_t - \beta_t \cdot \nabla \pi_{\omega_t}(\hat{s}_t, \hat{\mathbf{a}}_t), \quad (10)$$

where  $\beta_t$  is the learning rate that diminishes over time, i.e.,  $\lim_{t \rightarrow \infty} \beta_t = 0$ . The Robust-MA3C algorithm is summarized in **Algorithm 1**. An overview of the proposed method, i.e., Robust-MA3C is given in Fig. 2.

---

#### Algorithm 1 Robust-MA3C

---

- 1: **for**  $t = 1$  to maximum episode length **do**
  - 2:   Execute actions  $a = (a_1, \dots, a_N)$  and observe reward  $r$  and next states  $s' = (s'_1, \dots, s'_N)$ .
  - 3:   Store  $(s, a, r, s')$  in replay buffer  $\mathcal{D}$
  - 4:   **for** agent  $i$  to  $N$  **do**
  - 5:     Sample a random minibatch of samples  $(s, a, r, s')$  from  $\mathcal{D}$
  - 6:     Update the nature-agent's policies by using Eqn 10.
  - 7:     Update critic loss by using Eqn. 9.
  - 8:     Update actor loss by using Eqn. 8.
  - 9:   **end for**
  - 10: **end for**
- 

## V. EVALUATION

### A. Experiment Setup

The experiments reported here utilize a proprietary system-level network simulator. This simulator was created to simulate the behaviors of cellular communication networks. We use a network with 7 BSs to benchmark the proposed method, each of which supports 3 sectors. In each sector, there are 4 channels residing on 4 different carrier frequencies. (These 4 carrier frequencies are identical across different sectors and BSs.) The scenario is wrapped around at the edges. We emulate various traffic scenarios and days. The details of these simulation parameters can be seen in Table. I. Specifically, two types of UEs need to be balanced, with the two aforementioned load balancing mechanisms, i.e., AULB and IULB. These UEs are uniformly distributed geographically at initialization. The active UE movement follows a random walk process with an average speed of  $3m/s$ . The packet arrival follows a Poisson process with an average inter-arrival time of  $200ms$ .

To demonstrate the effectiveness of the proposed method, we conduct experiments in a scenario illustrated in Fig 1. In this scenario, there are 7 BSs, 6 arranged in a hexagonal ring

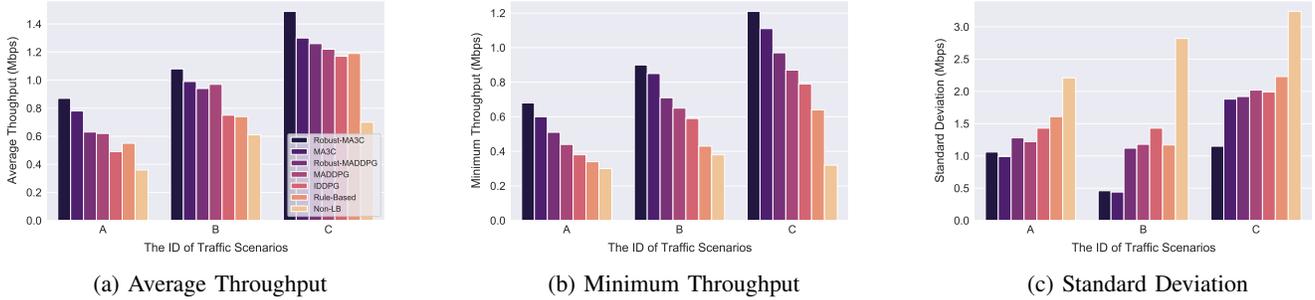


Fig. 3: System performance over various metrics in different traffic scenarios. In metric average throughput, the higher the better. In minimum throughput, the higher the better. In standard deviation, the lower the better.

Traffic ID	Day	UE	Active UE	Idle UE	Packet Size
A	1	40	12	28	0.41 Mbps
B	1	40	10	30	1.03 Mbps
C	1	40	7	33	0.89 Mbps
	2	40	18	22	1.18 Mbps
	3	40	19	21	1.11 Mbps

TABLE I: Parameters of different traffic scenarios and days.

around a central station. In this scenario, actions taken by one BS affect other BSs’ decisions. We control all of these 7 BSs to analyze the effectiveness of the proposed method.

### B. Methods Evaluated

We implement the following baseline methods. **Non-LB** evaluates the system performance without load balancing. **Rule-based** controls a BS based on prior knowledge and fixed control parameters. **Independent DDPG** trains each BS independently by utilizing DDPG policy learning algorithm [19]. **MADDPG** introduces centralized training decentralized execution framework to facilitate agents’ collaboration, which is one of the SOTA MARL methods [12]. **Robust-MADDPG** develops a Q-learning algorithm to find robust Nash equilibrium policies, which also solves the model uncertainty problem [8]. This is a follow-up work of MADDPG. **MA3C (ours)** the ablation study of proposed Robust-MA3C without a robust agent. Note that this approach does not guarantee to converge to the Nash equilibrium.

### C. Evaluation Results in Different Traffic Scenarios

We evaluate all the comparison schemes on different traffic scenarios, which are shown in Fig. 3 and showcase the evaluation results on the aforementioned metrics. From this figure, we can observe that the proposed MA3C and Robust-MA3C algorithm consistently outperforms other baselines in terms of various metrics. Compared to the third best baseline (i.e., Robust-MADDPG), MA3C increases the performance of average throughput, minimum throughput, and standard deviation by up to 12%, 31%, and 42%. This illustrates that, by utilizing the neighbor-aware attention mechanism, our proposed MA3C MARL algorithm enhances agents’ collaborations.

Furthermore, we visualize the effectiveness of combining robust agents by comparing Robust-MA3C to the third best baseline. The results show that Robust-MA3C improves the performance over corresponding metrics by up to 18%, 37%, and 45%. This verifies the effectiveness of finding Nash equilibrium and further advances the state-of-the-art.

### D. Evaluation Results in Different Days

Next, we further present results on different days of one traffic scenario, which are depicted in Fig. 4. We draw the same conclusion that the proposed two algorithms can improve the system performance overall metrics. Precisely, MA3C increases the performance of average throughput, minimum throughput, and standard deviation by up to 12%, 31%, and 42%. As for Robust-MA3C, it improves the performance over metrics by up to 19%, 37%, and 45%. These results further illustrate the effectiveness of the proposed algorithms.

### E. Evaluations of Individual BSs

To better understand the performance of each BS, we present a learning curve in Fig. 5. Note that we use the Robust-MA3C as the training algorithm as an example since MA3C illustrates the same phenomenon. In this figure, we observe that all BSs converge to the same range of rewards 1.25 ~ 1.5. This result illustrates that after training around 30k episodes, all BSs tend to achieve a Nash equilibrium and further illustrates the effectiveness of the proposed method. Among the learning curves, we can observe **an interesting phenomenon**: to reach a Nash equilibrium, some agents that achieve high reward at the early stage would decrease their performance for poorly performing BSs. For example, we find there is a big drop in the learning curve of BS#6. This BS#6 achieves a high reward at around 12k episodes. Then, the performance of this agent drastically drops to a low value. This happens because other agents cannot find optimal policies at this time. To help the poorly performing agents, BS#6 would like to explore some sub-optimal actions and sacrifice its performance to help the other agents. Finally, the received reward of BS#6 steadily increases after other agents converge. With the help of good agents, all agents collaborate and converge on the same range of rewards.

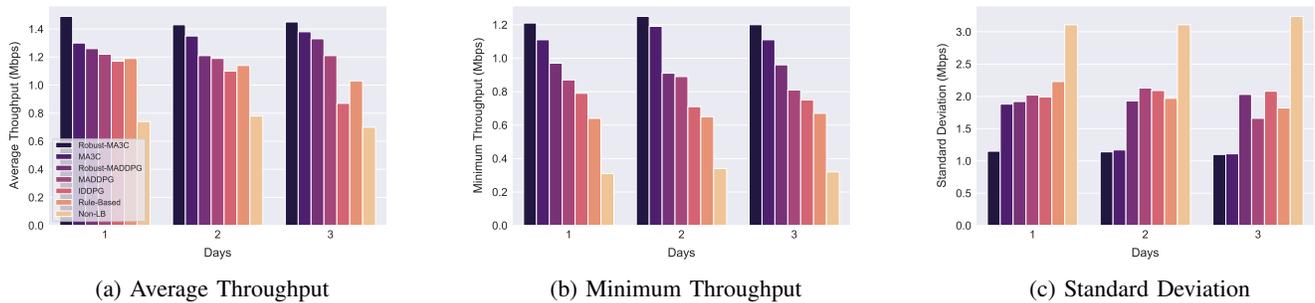


Fig. 4: System performance in different days.

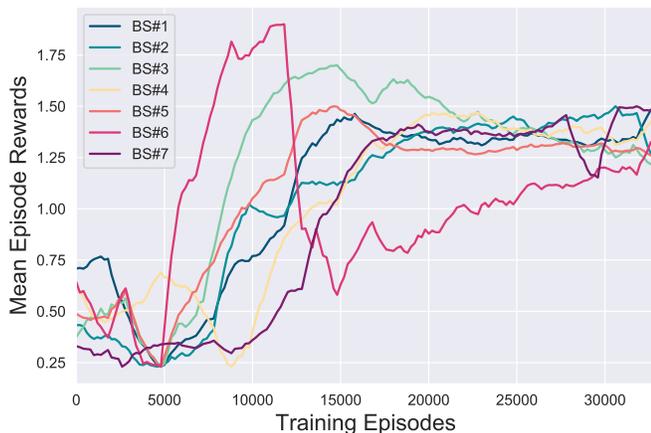


Fig. 5: The MA3C learning curves of seven BSs. The x-axis refers to the training episodes and the y-axis stands for the mean episode rewards that defined in Sec. III-C.

## VI. CONCLUSION

We have studied the multi-agent load balancing problem in cellular networks, where active UEs and idle UEs are migrated across BSs. A major challenge lies in how to enable collaboration among multiple BSs to achieve the joint Nash equilibrium policy for the agents. To address this challenge, we have proposed a novel algorithm, i.e., Robust-MA3C, a Robust-Multi-agent Attention Actor-Critic load balancing algorithm. Specifically, to provide a good migrating service, we utilize an attention mechanism to improve the performance of low-performance agents with the help of high-performance agents. In addition, we leverage the idea of nature-agent to achieve the Nash equilibrium. The evaluation results demonstrate that Robust-MA3C can achieve superior performance over the rule-based and other multi-agent reinforcement learning algorithms on three key network performance metrics. By using the simulation results, we have also demonstrated how the agents converge to Nash equilibrium.

## REFERENCES

- [1] L. Chiaraviglio, G. Bianchi, N. Blefari-Melazzi, and M. Fiore, "Will the proliferation of 5g base stations increase the radio-frequency "pollution"?", in *VTC Spring*, pp. 1–7, IEEE, 2020.
- [2] A. Feriani, D. Wu, Y. T. Xu, J. Li, S. Jang, E. Hossain, X. Liu, and D. Gregory, "Multi-objective load balancing for multi-band downlink cellular networks: A meta-reinforcement learning approach," *IEEE Journal on Selected Areas in Communications*, 2022, accepted, to appear.
- [3] Z. Yao, Z. Ding, and T. H. Clausen, "Reinforced workload distribution fairness," *CoRR*, vol. abs/2111.00008, 2021.
- [4] J. Kang, X. Chen, D. Wu, Y. T. Xu, X. Liu, G. Dudek, T. Lee, and I. Park, "Hierarchical policy learning for hybrid communication load balancing," in *ICC*, pp. 1–6, IEEE, 2021.
- [5] D. Wu, J. Kang, Y. T. Xu, H. Li, J. Li, X. Chen, D. Rivkin, M. Jenkin, T. Lee, I. Park, X. Liu, and G. Dudek, "Load balancing for communication networks via data-efficient deep reinforcement learning," in *GLOBECOM*, pp. 1–7, IEEE, 2021.
- [6] T. Mai, H. Yao, Z. Xiong, S. Guo, and D. T. Niyato, "Multi-agent actor-critic reinforcement learning based in-network load balance," in *GLOBECOM*, pp. 1–6, IEEE, 2020.
- [7] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *NIPS*, pp. 5998–6008, 2017.
- [8] K. Zhang, T. Sun, Y. Tao, S. Genc, S. Mallya, and T. Basar, "Robust multi-agent reinforcement learning with model uncertainty," in *NearIPS*, 2020.
- [9] Y. H. Kim, S. Han, C. H. Lyu, and H. Y. Youn, "An efficient dynamic load balancing scheme for multi-agent system reflecting agent workload," in *CSE (1)*, pp. 216–223, IEEE Computer Society, 2009.
- [10] J. Stender, S. Kaiser, and S. Albayrak, "Mobility-based runtime load balancing in multi-agent systems," in *SEKE*, pp. 688–696, 2006.
- [11] C. C. Myint and K. M. L. Tun, "A framework of using mobile agent to achieve efficient load balancing in cluster," in *6th Asia-Pacific Symposium on Information and Telecommunication Technologies*, pp. 66–70, IEEE, 2005.
- [12] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *NIPS*, pp. 6379–6390, 2017.
- [13] J. N. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," in *AAAI*, pp. 2974–2982, AAAI Press, 2018.
- [14] J. N. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson, "Learning to communicate with deep multi-agent reinforcement learning," in *NIPS*, pp. 2137–2145, 2016.
- [15] T. Chu, S. Chinchali, and S. Katti, "Multi-agent reinforcement learning for networked system control," in *ICLR*, OpenReview.net, 2020.
- [16] S. Sukhbaatar, A. Szlam, and R. Fergus, "Learning multiagent communication with backpropagation," in *NIPS*, pp. 2244–2252, 2016.
- [17] V. R. Konda and J. N. Tsitsiklis, "Actor-critic algorithms," in *NIPS*, pp. 1008–1014, The MIT Press, 1999.
- [18] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *NIPS*, pp. 1057–1063, The MIT Press, 1999.
- [19] 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 *ICLR (Poster)*, 2016.