

GAN-Powered Deep Distributional Reinforcement Learning for Resource Management in Network Slicing

Yuxiu Hua, Rongpeng Li[✉], Member, IEEE, Zhifeng Zhao, Xianfu Chen[✉], and Honggang Zhang[✉]

Abstract—Network slicing is a key technology in 5G communications system. Its purpose is to dynamically and efficiently allocate resources for diversified services with distinct requirements over a common underlying physical infrastructure. Therein, demand-aware resource allocation is of significant importance to network slicing. In this paper, we consider a scenario that contains several slices in a radio access network with base stations that share the same physical resources (e.g., bandwidth or slots). We leverage deep reinforcement learning (DRL) to solve this problem by considering the varying service demands as the environment *state* and the allocated resources as the environment *action*. In order to reduce the effects of the annoying randomness and noise embedded in the received service level agreement (SLA) satisfaction ratio (SSR) and spectrum efficiency (SE), we primarily propose generative adversarial network-powered deep distributional Q network (GAN-DDQN) to learn the action-value distribution driven by minimizing the discrepancy between the estimated action-value distribution and the target action-value distribution. We put forward a reward-clipping mechanism to stabilize GAN-DDQN training against the effects of widely-spanning utility values. Moreover, we further develop Dueling GAN-DDQN, which uses a specially designed dueling generator, to learn the action-value distribution by estimating the state-value distribution and the action advantage function. Finally, we verify the performance of the proposed GAN-DDQN and Dueling GAN-DDQN algorithms through extensive simulations.

Index Terms—Network slicing, deep reinforcement learning, distributional reinforcement learning, generative adversarial network, GAN, 5G.

Manuscript received June 20, 2019; revised October 15, 2019; accepted November 6, 2019. Date of publication December 12, 2019; date of current version February 19, 2020. This work was supported in part by the National Key R&D Program of China under Grant 2017YFB1301003, in part by the National Natural Science Foundation of China under Grant 61701439 and Grant 61731002, in part by the Zhejiang Key Research and Development Plan under Grant 2019C01002 and Grant 2019C03131, in part by the Zhejiang Lab under Grant 2019LC0AB01, in part by the Zhejiang Provincial Natural Science Foundation of China under Grant LY20F010016, and in part by the Fundamental Research Funds for the Central Universities under Grant 2019QN5010. This article was presented in part at the IEEE Globecom 2019 [1]. (*Corresponding author: Rongpeng Li*)

Y. Hua, R. Li, and H. Zhang are with the College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou 310027, China (e-mail: 21631087@zju.edu.cn; lirongpeng@zju.edu.cn; honggangzhang@zju.edu.cn).

Z. Zhao is with the Zhejiang Lab, Hangzhou 311121, China, and also with the College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou 310027, China (e-mail: zhaozf@zhejianglab.com).

X. Chen is with the VTT Technical Research Centre of Finland, 90570 Oulu, Finland (e-mail: xianfu.chen@vtt.fi).

Color versions of one or more of the figures in this article are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/JSAC.2019.2959185

I. INTRODUCTION

THE emerging fifth-generation (5G) mobile systems, armed with novel network architecture and emerging technologies, are expected to offer support for a plethora of network services with diverse performance requirements [2], [3]. Specifically, it is envisioned that 5G systems cater to a wide range of services differing in their requirements and types of devices, and going beyond the traditional human-type communications to include various kinds of machine-type communications [4]. According to ITU-R recommendations, it is a consensus that the key technologies of 5G wireless systems will spawn three generic application scenarios: enhanced mobile broadband (eMBB), massive machine-type communications (mMTC), and ultra-reliable and low-latency communications (URLLC) [5]. Specifically, (a) eMBB supports data-driven use cases requiring high data rates across a wide coverage area; (b) mMTC supports a very large number of devices in a broad area, which may only send data sporadically, such as Internet of Things (IoT) use cases; (c) URLLC supports strict requirements on latency and reliability for mission-critical communications, such as remote surgery, autonomous vehicles or Tactile Internet. Serving a diverse set of use cases over the same network will increase complexity, which must be managed to ensure acceptable service levels. As the services related to different use cases can have very different characteristics, they will impose varying requirements on the network resources. For example, considering mMTC devices for utility metering and parking sensors, they hardly move and do not require mobility management, i.e., the need to track location. On the other hand, sensors related to freight management would commonly move even across countries' borders and would require mobility management, including roaming agreements.

However, legacy mobile networks are mostly designed to provide services for mobile broadband consumers and merely consist of a few adjustable parameters like priority and quality of service (QoS) for dedicated services. Thus it is difficult for mobile operators to extend their networks into these emerging vertical services because of the different service requirements for network design and development. The concept of network slicing has recently been proposed to address this challenging problem; the physical and computational resources of the network can be sliced to meet the diverse needs of a range

of 5G users [2], [6], [7]. In this way, heterogeneous requirements can be served cost-effectively by the same physical infrastructure, since different network slice (NS) instances can be orchestrated and configured according to the specific requirements of the slice tenants.

As a non-nascent concept, network slicing can be traced back to the Infrastructure as a Service (IaaS) cloud computing model [8], whereby different tenants share computing, networking and storage resources to create different isolated fully-functional virtual networks on a common infrastructure. In the context of 5G and beyond, network functions virtualization (NFV) and software defined networking (SDN) technologies serve as a basis for the core network slicing by allowing both physical and virtual resources to be used to provide certain services, thus enabling 5G networks to deliver different kinds of services to various customers [9], [10]. On the other hand, the Next Generation Mobile Networks (NGMN) alliance puts forward an evolved end-to-end network slicing idea while the Third-Generation Partnership Project (3GPP) also suggests that radio access network (RAN) should not be excluded to “design specific functionality to support multiple slices or even partition of resources for different network slices” [11], [12].

However, in order to provide better-performing and cost-efficient services, RAN slicing involves more challenging technical issues for the realtime resource management on existing slices, since (a) for RANs, spectrum is a scarce resource and it is essential to guarantee the spectrum efficiency (SE) [9]; (b) the service level agreements (SLAs) with slice tenants usually impose stringent requirements; and (c) the actual demand of each slice heavily depends on the request patterns of mobile users (MUs) [13]. Therefore, the classical dedicated resource allocation fails to address these problems simultaneously [11]. Instead, it is necessary to intelligently allocate the radio resources (e.g., bandwidth or slots) to slices according to the dynamics of service requests from mobile users coherently [14] with the goal of meeting SLA requirements in each slice, but at the cost of acceptable SE. In this regard, there have been extensive efforts [13], [15]–[20], [43]. [15] proposed an online genetic slicing strategy optimizer for inter-slice resource management. However, [15] did not consider the explicit relationship between the required resource and SLA on a slice, as one slice might require more resources given its more stringent SLA. Reference [16] considered the problem of different types of resources (bandwidth, caching, backhaul capacities) being allocated to NS tenants based on user demands. The authors proposed mathematical solutions, but the optimization problem would become intractable when the simulation parameters are scaled up (e.g., increasing the number of NSs or the shareable resources). Reference [17] mathematically analyzed the joint optimization problem of access control and bandwidth allocation in the multi-base station (BS) multi-NS scenario. However, the solutions therein are based on the assumption that different users have the same fixed demand rate, which is a condition unlikely to be found in practice. From the perspective of bandwidth usage-based pricing (UBP), [20] used game theory to analyze the relationship between Internet service providers (ISPs) and users, thereby improving the profit of ISPs and solving

the peak-time congestion problem. However, the timeslot for bandwidth allocation in [20] was 1 hour, which is unrealistic in a situation where the number of demands changes drastically in a short period.

In order to address the demand-aware resource allocation problem, one potential solution is reinforcement learning (RL). RL is an important type of machine learning where an agent learns how to perform optimal actions in an environment from observing state transitions and obtaining feedback (rewards/costs). In RL, the action value, $Q(s, a)$, describes the expected return, or the discounted sum of rewards, when performing action a in state s . Usually, the Q value can be estimated by classic value-based methods such as SARSA [21] and Q-learning [22] based on the Bellman equation. Reference [23] used deep neural networks to approximate the Q function, namely deep Q network (DQN), which demonstrated human-like performance on simple computer games and inspired a research wave of deep reinforcement learning (DRL). Besides modeling $Q(s, a)$, [24] showed that we could learn the distribution of $Q(s, a)$ by a distributional analogue of Bellman equation; this approach improved the estimation of action values in an inherently randomness environment. Specifically, [24] proposed C51 algorithm to minimize the Kullback–Leibler (KL) divergence between the approximated Q distribution and the target Q distribution calculated by the distributional Bellman optimality operator. Inspired by the theory of quantile regression [25], [26] proposed the quantile regression DQN (QR-DQN) and thus successfully performed distributional RL over the Wasserstein metric, leading to the state-of-the-art performance. Reference [27] extended QR-DQN from learning a discrete set of quantiles to learning the full quantile function and put forward the implicit Q network (IQN). Given the success of replacing $Q(s, a)$ by its distribution in [24], [26], [27] as well as the reputation of generative adversarial network (GAN) for approximating distributions [28], it naturally raises a question of whether GAN is viable for approximating the action-value distribution and thus improving distributional RL.

In the field of communications and networking, DRL has triggered tremendous research attention to solving resource allocation issues in some specific fields like power control [29], green communications [30], cloud RANs [31], mobile edge computing and caching [32]. Given the challenging technical issues in resource management on existing NSs, the previous work in [13] leveraged DQN to find the optimal resource allocation policy and investigated its performance. However, the method proposed in [13] did not consider the effects of random noise on the calculation of SE and SLA satisfaction ratio (SSR). To mitigate the potential risk of incorrectly estimating the action value due to the randomness in SE and SSR, we intend to introduce the distributional RL to estimate the action-value distribution, thus avoiding the action-value overestimation or underestimation issue that plagues many traditional value-based RL algorithms (e.g., DQN). Meanwhile, the cutting-edge performance of Wasserstein generative adversarial network with gradient penalty (WGAN-GP) in the distribution approximation suggests to us that we might use it to learn the

action-value distribution. To this end, we propose a new approach, the GAN-powered deep distributional Q network (GAN-DDQN), based on distributional RL and WGAN-GP, to realize dynamic and efficient resource allocation per slice.

The main contributions of this paper are as follows:

- To find the optimal resource allocation policy under the uncertainty of slice service demands, we design the GAN-DDQN algorithm, where the generator network outputs a fixed number of particles that try to match the action-value distribution for each action. Such a design in GAN-DDQN can mitigate the effects of learning from a nonstationary environment and is significantly different from the concurrent yet independent works [33], [34].¹
- We demonstrate that the widely-spanning system utility values could destabilize GAN-DDQN's training process, and correspondingly design a reward-clipping mechanism to reduce this negative impact. Specifically, we clip the system utility values to some constant values according to a straightforward rule with several heuristics-guided adjustable thresholds, and then use these constants as the final rewards in RL.
- GAN-DDQN suffers from the challenge that only a small part of the generator output is included in the calculation of the loss function during the training. To compensate for this, we further propose Dueling GAN-DDQN, which is a special solution derived from Dueling DQN [35] and the discrete normalized advantage functions (DNAF) algorithm [36]. Dueling GAN-DDQN separates the state-value distribution from the action-value distribution and combines the action advantage function to obtain the action values. In addition, we elaborate on twofold loss functions that further take advantage of the temporal difference (TD) error information to achieve performance gains. The introduction of dueling networks to GAN-DDQN makes the work described in this paper significantly different from our previous work presented in IEEE Globecom 2019 [1].
- Finally, we perform extensive simulations to demonstrate the superior efficiency of the proposed solutions over the classical methods, such as DQN, and provide insightful numerical results for the implementation details.

The remainder of the paper is organized as follows: Section II talks about some necessary mathematical backgrounds and formulates the system model. Section III gives the details of the GAN-DDQN, while Section IV presents the detailed simulation results. Finally, Section V summarizes the paper and offers prospects.

II. PRELIMINARIES AND SYSTEM MODEL

A. Preliminaries

Table I lists the important notations used in this paper. An agent tries to find the optimal behavior in a given setting

¹ [33] used a generator network that directly outputs action values and did not show any significant improvement of GAN Q-learning over conventional DRL methods. [34] used the policy iteration method [37] to loop through a two-step procedure for the value estimation and policy improvement, where GAN was only used to estimate the action-value distribution. Besides, [34] exploited a totally different framework without the target generator, which is a key component for GAN-DDQN.

TABLE I
NOTATIONS USED IN THIS PAPER

Notation	Definition
\mathcal{S}	State space
\mathcal{A}	Action space
P	Transition probability
V	State-value function
Q	Action-value function
Z_v	Random variable to statistically model the state values
Z_q	Random variable to statistically model the action values
\mathcal{T}^*	Bellman optimality operator
s, s'	States
a	An action
r	A reward
S_t	State at time t
A_t	Action at time t
R_t	Reward at time t
γ	Discount factor
π	Policy
J	System utility
α	Weight of the SE
β	Weight of the SSR
τ	Quantile samples
λ	Gradient penalty coefficient
n_{critic}	The number of discriminator updates per training iteration

↑ 這段是對方程 B-E, Greedy Optimal Policy, F.A.

through interaction with the environment, which can be treated as solving an RL problem. This interactive process can be modeled as a Markov Decision Process $(\mathcal{S}, \mathcal{A}, R, P, \gamma)$, where \mathcal{S} and \mathcal{A} denote the state and action spaces, R is the reward, $P(\cdot|s, a)$ is the transition probability, and $\gamma \in (0, 1]$ is a discount factor. A policy $\pi(\cdot|s)$ maps a state to a distribution over actions. The state-value function of a state s under a policy $\pi(\cdot|s)$, denoted $V^\pi(s)$, is the expected return when starting in s and following π thereafter. Similarly, we define the value of taking action a in state s under the policy π , denoted $Q^\pi(s, a)$, as the expected return starting from s , taking the action a , and thereafter following policy π . Mathematically, the state-value function is

$$V^\pi(s) = \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{\infty} \gamma^t R_t | S_0 = s \right], \quad (1)$$

and the action-value function is

$$Q^\pi(s, a) = \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{\infty} \gamma^t R_t | S_0 = s, A_0 = a \right], \quad (2)$$

where \mathbb{E} denotes the expectation. The relationship between the value of a state and the values of its successor states is expressed by the Bellman equation for V^π

$$V^\pi(s) = \mathbb{E}_{\pi, P} [R + \gamma V^\pi(s')]. \quad (3)$$

Similarly, the Bellman equation for Q^π is

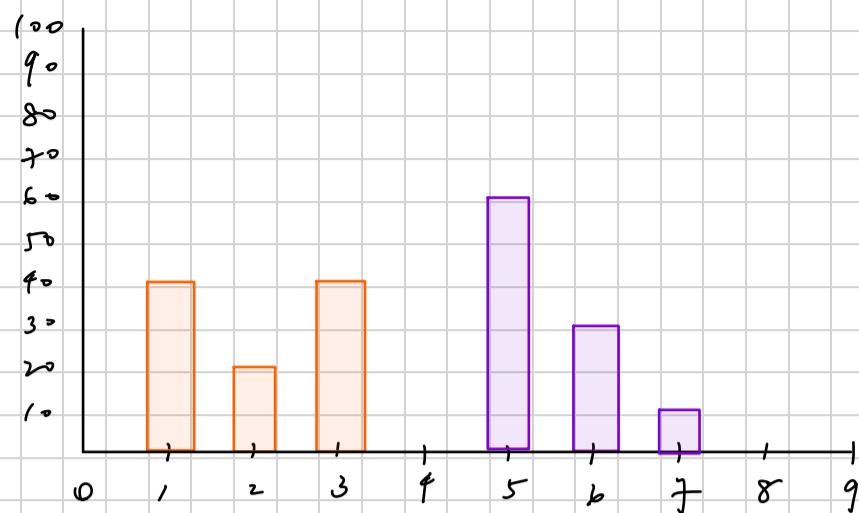
$$Q^\pi(s, a) = \mathbb{E}_{\pi, P} [R + \gamma Q^\pi(s', a')], \quad (4)$$

where s' and a' can be derived from the transition probability $P(\cdot|s, a)$ and a policy $\pi(\cdot|s')$, respectively.

The goal of RL is to find the optimal policy which yields the maximum $Q(s, a)$ for all s and a . Let $\pi^* = \arg \max_\pi Q^\pi(s, a)$ be the optimal policy and let $Q^*(s, a)$ be the corresponding

1-Wasserstein Distance

Q



$$X \quad \begin{cases} P(X=0) = 0.4 \\ P(X=2) = 0.2 \\ P(X=3) = 0.4 \end{cases}$$

$$Y \quad \begin{cases} P(Y=5) = 0.6 \\ P(Y=6) = 0.3 \\ P(Y=7) = 0.1 \end{cases}$$

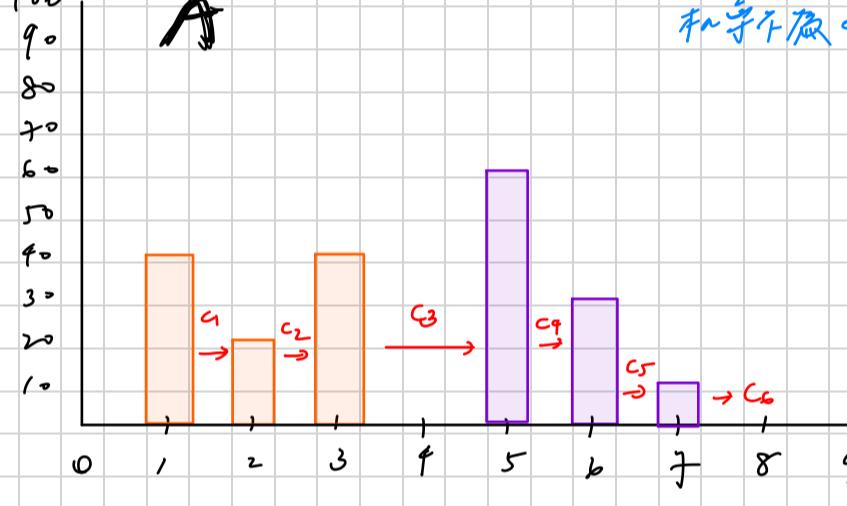
$$\rho\text{-Wasserstein Distance} \quad W_p(X, Y) = \left(\int_0^1 |F_X^{-1}(w) - F_Y^{-1}(w)|^p dw \right)^{1/p}$$

\Rightarrow Demo \rightarrow 1-Wasserstein Distance by PMF (Discrete)

$$W_1(X, Y) = \sum_{x_i=1}^n |F_p(x_i) - F_q(x_i)| \cdot \Delta x$$

質量 拼文連距離 $= x_{i+1} - x_i$

* 這裡用 COF 而非 Inverse COF；離散時它們理論上等價。



和字不麻。百又莫。

$$\textcircled{1} \quad \text{Support} = \{1, 2, 3, 5, 6, 7\}$$

\textcircled{2} COF of each x of X & Y :

$$F_X(x) = P(X \leq x)$$

$$F_Y(x) = P(Y \leq x)$$

$$x \quad F_X(x)$$

$$F_Y(x)$$

$$1 \quad 0.4$$

$$0$$

$$2 \quad 0.6 (0.2 + 0.4)$$

$$0$$

$$3 \quad 1 (0.2 + 0.4 + 0.4)$$

$$0$$

$$5 \quad 1$$

$$0.6$$

$$6 \quad 1$$

$$0.9$$

$$7 \quad 1$$

$$1$$

前面推來的 0.9 跟現有 0.2

\textcircled{3}

$$x \quad |F_X(x) - F_Y(x)| \quad \Delta x$$

C1

$$1 \quad 0.4 - 0 = 0.4$$

1

X 能行右推而上 跟下一致距離

C2

$$2 \quad 0.6 - 0 = 0.6$$

1

前面推來的 0.9 跟現有 0.2

C3

$$3 \quad 1 - 0 = 1$$

2

在 2, 2, 3 中 X 能行左推 跟下一致距離

C4

$$5 \quad 1 - 0.6 = 0.4$$

1

在 5 的 2 位罝用 0.6 去堆 Y3, 剩 0.4 能行左推

C5

$$6 \quad 1 - 0.9 = 0.1$$

1

C6

$$7 \quad 1 - 1 = 0$$

1

$$\text{Sum} = 0.4 + 0.6 + 2 + 0.4 + 0.1 = 3.5$$

Def T^* & D.B.E. Operator:

Contraction Mapping

$$T^\pi Z_q(s, a) \stackrel{D}{=} R + \gamma \cdot Z_q(s', a')$$

$$\|f(x_1) - f(x_2)\| \leq \Gamma \cdot \|x_1 - x_2\| \quad Z_q^\pi(s, a) \stackrel{D}{=} R + \gamma Z_q^\pi(s', a'),$$

(A)

Bellemare et al. proved that (the distributional Bellman equation is a contraction in p -Wasserstein distance) but the (distributional Bellman optimality operator is not necessarily a contraction, which provides a guideline for metric selection)

$$T^* Z_q(s, a) \stackrel{D}{=} R + \gamma Z_q \left(s', \arg \max_{a' \in \mathcal{A}} \mathbb{E}[Z_q(s', a')] \right)$$

① $\mathbb{E} W_p(T^\pi z_1, T^\pi z_2) \leq \Gamma \cdot W_p(z_1, z_2)$,

代表 $T^\pi z_1$ & $T^\pi z_2$ 有 p -Wasserstein Distance

即 z_1 & z_2 有 p -W.D. 且 Γ 越小。

② $W_p(T^\pi z_1, T^\pi z_2) \leq \Gamma \cdot W_p(z_1, z_2)$

若在 KL-Divergence 的意义下，D.B.E. 就不会收敛。

故在设计 Loss function 时会用 p -Wasserstein Distance。

EP model: $\hat{z}_\theta(s, a)$

理論上收斂率高

Target: $z_{\text{target}} = r + \gamma \cdot \underline{z}(s', a')$

$$\text{Loss: } L(\theta) = W_p(\hat{z}_\theta(s, a), T^\pi z(s, a))$$

但前面不是說要最小化

$$\sup_{s, a} \text{dist}(T^* Z_q(s, a), Z_q(s, a)),$$

但後面又講到 DBDO(T^*)

是這樣子，雖然沒有理論上收斂率高，

但我們仍能設計一個合理的 Loss function (by p-Was.)

畢竟本來就很難滿足理論的 R.L.

(B) A. 前面那段說：

The distributional Bellman optimality operator is not necessarily a contraction, which provides a guideline for metric selection.

✓ 這是 Bellemare et al. 2017 的原始分析 (C51 論文) :

- 在某些距離空間 (如 KL 或 TV 距離) , DBOO 不是 contraction
- 甚至連在 1-Wasserstein distance 上也不是, 除非你有額外假設

→ 改 C51 使用 KL-Divergence

} 改善收斂性及
穩定性

B. 現在你這段說：

DBOO is a strict contraction in p-Wasserstein distance

✓ 這來自的是 QR-DQN (Dabney et al., 2018) 的分析, 它針對的是:

分布是透過 quantile 表示 (support 是固定位置) 的特定 case,
在這種架構下, DBOO 的確可以證明為 strict contraction in W_p 。

→ 改 QR-DQN 使用 p-Wasserstein

核心關鍵點：

條件	是否是 contraction
一般隨機變數分布 (無固定 support)	✗ DBOO 不是 contraction in W_1
分布以 quantile 表示 (QR-DQN 設計)	✓ DBOO 是 strict contraction in W_p

也就說：不是所有 case 都有 contraction，但某些結構設計下是可以證明 contraction 的。

更具體地說：

- 若你令 return distribution 表示為：

$$Z(s, a) = \sum_{i=1}^N \delta_{z_i} \cdot \frac{1}{N}$$

並固定 support 的位置 (像 QR-DQN) , 就可以證明 DBOO 是 contraction in W_1

- 但若分布是任意形式，或是 support 不是固定的，則 contraction 保證會失效

總結：

1. 固定 Quantile levels, 2. 有限支撑集

廣義而言, DBOO 不保證 Contraction.

但在 QR-DQN 這種結構中 Contraction.

Strict Contraction $\| f(x_1) - f(x_2) \| < \| x_1 - x_2 \|$

→ in p-Wasserstein $W_p(f(x_1), f(x_2)) < W_p(x_1, x_2)$

💡 那直接最小化 Wasserstein 可以嗎？

✗ 在原始形式下不太可行，原因如下：

1. Wasserstein 距離需要排序或 inverse CDF，這些是不可微操作（梯度無法傳遞）
2. 你需要知道兩個分布的完整形式或 samples，而不是單一點估計
3. 無法直接在神經網路參數上定義 gradient

⌚ 世界上主要有兩種方法在逼近 Wasserstein Distance：

方法類型	代表技術	說明
✓ 間接法（可微）	Quantile Regression (QR-DQN)	針對一維機率分布，利用 quantile loss 近似 W_1 ，可以直接反向傳播
✓ 對偶法 (GAN 用)	Wasserstein GAN (WGAN)	使用 Kantorovich–Rubinstein 對偶定義，導入 critic 來估計 W_1 ，間接逼近

✗ WGAN 怎麼讓 f 是 1-Lipschitz？

常見的方法有：

方法	說明
Weight Clipping (傳統 WGAN 作法)	將每層神經網路的權重限制在一個區間（如 $[-0.01, 0.01]$ ）
Gradient Penalty (WGAN-GP)	加一項 loss：懲罰梯度 norm 偏離 1 的程度
Spectral Normalization	對每層權重做奇異值正規化，保證 Lipschitz 常數不超過 1

傳統 WGAN 的工作方法為使用 weight Clipping 作為來 + 加上 1-Lipschitz function

→ 但這種工作方法太圓滑，容易令訓練不穩定

WGAN-GP 改用 Gradient Penalty 的方式來改善。

→ 不再硬性 Clipping，而是以 G.P. 來軟性約束，以穩定訓練。

C51 → Q值切成51等分，minimize KL-Divergence

QR-DQN → 把Q值切成若干 Quantile，minimize 1-Wasserstein by Quantile Regression Loss

* temporal consistency: 从高時刻到低時刻保持一致。文中說的是不論時間順序， Q^* 的值會保持是正確的。

* Greedy Optimal Policy: 當了 state 選最大 Q 值的 Action，就能找到 it。* robustness: 一介於模型在那種情況下，環境變化時仍能維持正確的表現。

$\overline{Q}_{\text{optimal}} \text{ Q-value}$

action-value function. $Q^*(s, a)$ satisfies the following Bellman optimality equation

$$Q^*(s, a) = \mathbb{E}_{\pi^*, P} \left[R + \gamma \max_{a' \in \mathcal{A}} Q^*(s', a') \right]. \quad (5)$$

Eq. (5) illustrates the temporal consistency of the action-value function, which allows for the design of learning algorithms. Define the Bellman optimality operator T^* as

$$T^* Q(s, a) = \mathbb{E}_{\pi, P} \left[R + \gamma \max_{a' \in \mathcal{A}} Q(s', a') \right]. \quad (6)$$

When $\gamma \in (0, 1)$, starting from any $Q_t(s, a)$, iteratively applying the operator $Q_{t+1}(s, a) \leftarrow T^* Q_t(s, a)$ leads to convergence $Q_t(s, a) \rightarrow Q^*(s, a)$ as $t \rightarrow \infty$ [37].

In high dimensional cases, it is critical to use function approximation as a compact representation of action values. Let $Q_\theta(s, a)$ denote a function with parameter θ that approximates a table of action values with entry (s, a) . The optimization aim is to find θ such that $Q_\theta(s, a) \approx Q^*(s, a)$, and the optimal solution can be found by iteratively leveraging the Bellman optimality operator T^* . In other words, the optimal parameter θ can be approached by minimizing the squared TD error

$$\zeta^2 = \left[r + \gamma \max_{a' \in \mathcal{A}} Q_\theta(s', a') - Q_\theta(s, a) \right]^2 \quad (7)$$

over samples (s, a, r, s') , which are randomly selected from a replay buffer [38] that stores transitions which record the interaction between an agent and the environment when following the policy driven by Q_θ . In cases where $Q_\theta(s, a)$ is linear, the iterative process to find $Q^*(s, a)$ can be shown to converge [39]. However, in cases where $Q_\theta(s, a)$ is nonlinear (e.g., a neural network), $Q_\theta(s, a)$ becomes more expressive at the cost of no convergence guarantee. A number of DRL algorithms are designed following the above formulation, such as DQN [23] and Dueling DQN [35].

1) Distributional Reinforcement Learning: The main idea of distributional RL [24] is to work directly with the distribution of returns rather than their expectation (i.e., Q^π), so as to increase robustness to hyperparameter variation and environment noise [40]. Let the random variable $Z_q^\pi(s, a)$ be the return obtained by following a policy π to perform action a from the state s . Notably, the value of $Z_q^\pi(s, a)$ varies due to the unexpected randomness in the environment. Then we have

$$Q^\pi(s, a) = \mathbb{E}[Z_q^\pi(s, a)], \quad (8)$$

and an analogous distributional Bellman equation, that is,

$$Z_q^\pi(s, a) \stackrel{D}{=} R + \gamma Z_q^\pi(s', a'), \quad (9)$$

where $A \stackrel{D}{=} B$ denotes that random variable A has the same probability law as B . Therefore, a distributional Bellman optimality operator T^* can be defined by

$$T^* Z_q(s, a) \stackrel{D}{=} R + \gamma Z_q \left(s', \arg \max_{a' \in \mathcal{A}} \mathbb{E}[Z_q(s', a')] \right). \quad (10)$$

In traditional RL algorithms, we seek the optimal Q function approximator by minimizing a scalar value ζ^2 in Eq. (7).

Compare with (9)

$$Q^\pi(s, a) = \mathbb{E}_{\pi, P} [R + \gamma Q^\pi(s', a')],$$

Compare with (10)

$$T^* Q(s, a) = \mathbb{E}_{\pi, P} [R + \gamma \max_{a' \in \mathcal{A}} Q(s', a')].$$

做了 map 在進行如何各個特質是獨立的，這部分會用 $\sup_{s, a}$ 來定義 Objective function，

用來做跟 π 有關的。而且又可以得到

In distributional RL, our objective is to minimize a statistical distance: 意思是“最大誤差的行動策略”，即最壞情況下的優勢。

$$\sup_{s, a} \text{dist} \left(\frac{T^* Z_q(s, a)}{\text{Optimal Q.P.O.}}, \frac{Z_q(s, a)}{\text{Current Q.P.D.}} \right), \quad (11)$$

這裡的 dist 距離是 p -D.

因為這距離是到 p -D.

所以這距離是 p -Divergence.

在上邊，距離是到 p -D.

WGAN-GP Obj. function 解析

G: 產真圖片, D: 舉圖片打偽率, 越高越好

$$\min_G \max_{D \in \mathcal{D}} \frac{\mathbb{E}_{x \sim p_{\text{data}}} [D(x)] - \mathbb{E}_{z \sim p_z(z)} [D(G(z))]}{②} + p(\lambda),$$

- ① $\mathbb{E}_{x \sim p_{\text{data}}} [D(x)]$ p_{data} : 真美的圖片集之機率分布
 $x \sim p_{\text{data}} \rightarrow x$, sample from p_{data} , 遵循 p_{data} 之機率分布
 $D(x) \rightarrow x$ 故真美度分布
 → 代表 D 舉真美圖片打偽率美度分布期望值。
- ② $\mathbb{E}_{z \sim p_z(z)} [D(G(z))]$ $p_z(z) \rightarrow z$ 常分布
 $z \sim p_z(z) \rightarrow z$, sample from $p_z(z)$, 遵循 p_z 之機率分布
 $G(z) \rightarrow G$ 根據 z 產出的圖片
 $D(G(z)) \rightarrow D$ 舉 $G(z)$ 打偽率美度分布
 → 代表 D 舉 G 產生的假美圖所打偽率美度分布期望值

回到，該 loss function 結合 G & D 的自身乘成一式子：

$$① D : \max_D \mathbb{E}_{x \sim p_{\text{data}}} [D(x)] - \mathbb{E}_{z \sim p_z(z)} [D(G(z))]$$

→ 想最大化真圖及假圖偽差異  產真對抗訓練

$$② G : \min_G \mathbb{E}_{x \sim p_{\text{data}}} [D(x)] - \mathbb{E}_{z \sim p_z(z)} [D(G(z))]$$

→ 想最小化真圖及假圖偽差異

當二者達到 Nash 平衡時， $G(z)$ 會和 P_{data} 重合，代表 GAN 成功學習，即 G 成功達到真美的資料分布。 G 不能再更真，D 不能再偽 (ex. D 雖以 0.5)

★ 理論上，訓練不會真的達到完美 Nash，但只要 G 能產生更美的樣本即能成功訓練 GAN 成功訓練出來。★ D 的學習速度會“不肖快”。

* Gradient Penalty (GP)

為了使得 D 滿足 1-Lipschitz 的要求。

基本想法就是將 D 的梯度限制其範疇在 ≈ 1

→ 用來檢查 1-Lipschitz 是否成立的檢查式。

① $\hat{x} = \varepsilon \underline{\text{real}} + (1 - \varepsilon) \underline{\text{fake}}$, $\varepsilon \sim U(0, 1)$.
→ 決定 real / fake 的比例比值

② $p(\lambda) = \frac{\lambda}{2} \left(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1 \right)^2$
hyperparameter
→ L_2 norm, 表長度 or 范數
→ 反正坡偏離 (1-Lipschitz 等於 1)
→ 檢查是否在 D 附近做變化等
(1-Lipschitz 限制莫大於 1)

這種不直接給閾值的 weight clipping 的操作方法，能更穩定，較不會有梯度消失的問題。

* System Model

state: $[d_1, d_2, d_3]$



欲使用 2nd N.S. 传输的数据

packets 数

Action Space: $[w_1, w_2, w_3]$ 在此中，每一时隙都有三个频段可供资源分配。

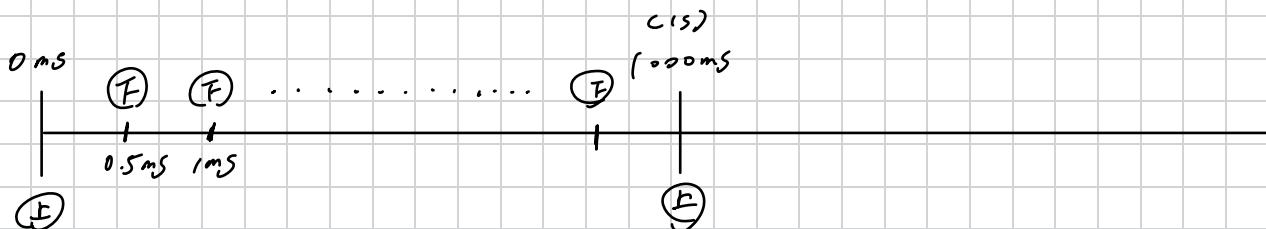
→ Total Bandwidth = 10 MHz

Bandwidth Allocation Resolution = 1 MHz

(Granularity)

$$\begin{array}{l} \left| \begin{array}{l} w_1 + w_2 + w_3 = 10 \\ w_i \geq 1 \end{array} \right. \rightarrow \text{所有时隙集合，其 } C_{\frac{J+3-1}{2}}^q = C_2^9 = 36 \\ \left(\begin{array}{l} k_1 + k_2 + k_3 = 7 \\ k_n \geq 0 \end{array} \right) \quad \text{Action Space} \end{array}$$

* Time Sequence



(上): 将 Bandwidth 分给 N.S. by GAN-DQCN: 每 15 毫秒行一次

(下): 将 N.S. 分到的 Band. 分给 active UE, by RR: 每 0.5 毫秒行一次

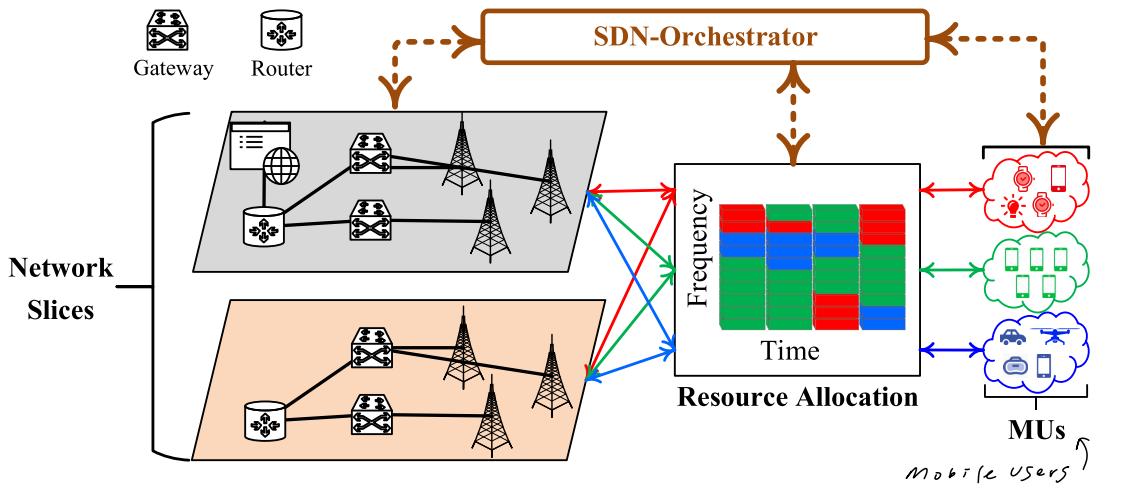


Fig. 1. The considered scenario showing uplink and downlink transmissions on different NSs.

把 NN 單位的反梯度限制在一區間內。 Lipschitz Function 是單調增強於 1~1, 但 GP 是讓它靠近 1

constraint by (clipping the weights of the discriminator to be within a certain range) that is governed by a hyperparameter. Furthermore, [42] proposed WGAN-GP and adopted gradient penalty to enforce the 1-Lipschitz constraint instead of simply clipping weights. Its optimization objective is formulated as follow:

$$\min_{G} \max_{D \in \mathcal{D}} \mathbb{E}_{x \sim p_{\text{data}}} [D(x)] - \mathbb{E}_{z \sim p_z(z)} [D(G(z))] + p(\lambda), \quad (13)$$

where \mathcal{D} denotes the set of 1-Lipschitz functions, x denotes the samples from real data, z denotes the samples from a random distribution, and $p(\lambda) = \frac{\lambda}{2} (\|\nabla_x D(\hat{x})\|_2 - 1)^2$, $\hat{x} = \varepsilon x + (1 - \varepsilon)G(z)$, $\varepsilon \sim U(0, 1)$. Gradient penalty increases the computational complexity but it does make WGAN-GP perform much better than previous GANs.

$$\text{梯度懲罰 } p(\lambda) = \sum_i [\max(0, \| \nabla_x D(\hat{x}_i) \|_2 - 1)]^2$$

B. System Model

Fig. 1 illustrates the SDN-based system model for dynamic allocation of wireless bandwidth in the RAN scenario with uplink and downlink transmissions. We consider the downlink case in this paper. In this respect, [43] built an SDN-based C-RAN testbed that realizes the dynamic allocation of radio resources (e.g., wireless bandwidth) by using frequency division duplex scheme, and the demonstration was presented in [44]. Under the framework of hierarchical network slicing, we consider a RAN scenario with a single BS, where a set \mathcal{N} of NSs share the aggregated bandwidth W .³ The bandwidth is allocated to each NS according to the number of demands for the corresponding type of service. For an NS, say NS n , it provides a single service for a set of users \mathcal{U}_n . We consider a timeslot model where the slicing decision is updated according to the demand of users periodically (e.g., 1 second). In one timeslot, the number of demands NS n receives is denoted as d_n , which partially determines the wireless bandwidth that the BS allocates to this NS, denoted as w_n .

³In fact, such a bandwidth allocation could be realized by physical multiplexing methods [45]. Meanwhile, the shared resources could be temporal slots as well. However, for simplicity of representation, we take the bandwidth allocation problem as an example.

The objective of our work is to find an optimal bandwidth-allocation solution that maximizes the system utility, denoted by J , which can be described by the weighted sum of SE and SSR. We now study the two sub-objectives, respectively. Let r_{un} be the downlink rate of user u_n served by NS n , which is, for simplicity, defined by Shannon theory as follows $w_n \rightarrow u_n$ 在 1 timeslot 累積 B.W.

$$r_{un} = w_n \log(1 + \text{SNR}_{un}), \quad \forall u_n \in \mathcal{U}_n, \quad (14)$$

where SNR_{un} is the signal-to-noise-ratio between user u_n and the BS. SNR_{un} can be given as

$$\text{SNR}_{un} = \frac{g_{un} P_{un}}{N_0 w_n}, \quad \begin{matrix} \xrightarrow{\text{Signal}} & \text{Unit Power Watt} \\ \xrightarrow{\text{Noise}} & \text{large-scale} \end{matrix} \quad (15)$$

where g_{un} is the average channel gain that captures (path loss and shadowing) from the BS to the user u_n , P_{un} is the transmission power, and N_0 is the single-side noise spectral density. Given the transmission rate, SE can be defined as follows

$$\text{SE} = \frac{\sum_{n \in \mathcal{N}} \sum_{u_n \in \mathcal{U}_n} r_{un}}{W}, \quad \begin{matrix} \xrightarrow{\text{每個子網有多少大數量的傳輸速率}} & \text{in every timeslot} \\ \xrightarrow{\text{累積}} & \end{matrix} \quad (16)$$

On the other hand, SSR of NS n is obtained by dividing the number of successfully transmitted packets by the total number of arrived packets on NS n . Before formulating this problem, we define \mathcal{Q}_{un} as the set of packets sent from the BS to user u_n , determined by the actual traffic demand patterns, and define a binary variable $x_{qu_n} \in \{0, 1\}$, where $x_{qu_n} = 1$ indicates that the packet $qu_n \in \mathcal{Q}_{un}$ is successfully received by user u_n , i.e., the downlink data rate r_{un} and the latency l_{qu_n} are simultaneously satisfied. Therefore, $x_{qu_n} = 1$ if and only if $r_{un} \geq \bar{r}_n$ and $l_{qu_n} \leq \bar{l}_n$, where l_{qu_n} denotes the latency that takes account of both queuing delay and transmission delay. \bar{r}_n and \bar{l}_n are the predetermined rate and latency values according to the SLA for service type n . We can formulate the SSR for NS n as: Success Satisfaction Ratio.

$$\text{SSR}_n = \frac{\sum_{u_n \in \mathcal{U}_n} \sum_{qu_n \in \mathcal{Q}_{un}} x_{qu_n}}{\sum_{u_n \in \mathcal{U}_n} |\mathcal{Q}_{un}|}, \quad \begin{matrix} \xrightarrow{\text{達成SLA的Packet}} & \text{達成SLA的Packet} \\ \xrightarrow{\text{總}} & \text{總 Packet} \end{matrix} \quad (17)$$

where $|\mathcal{Q}_{un}|$ denotes the number of packets sent from the BS to user u_n . * 在 NS 中的 SLA

VOLTE: rate $> 51 \text{ kbps}$, latency $< 10 \text{ ms}$
 eMBB: rate $> 100 \text{ Mbps}$, latency $< 10 \text{ ms}$
 uRLLC: rate $> 10 \text{ Mbps}$, latency $< 1 \text{ ms}$

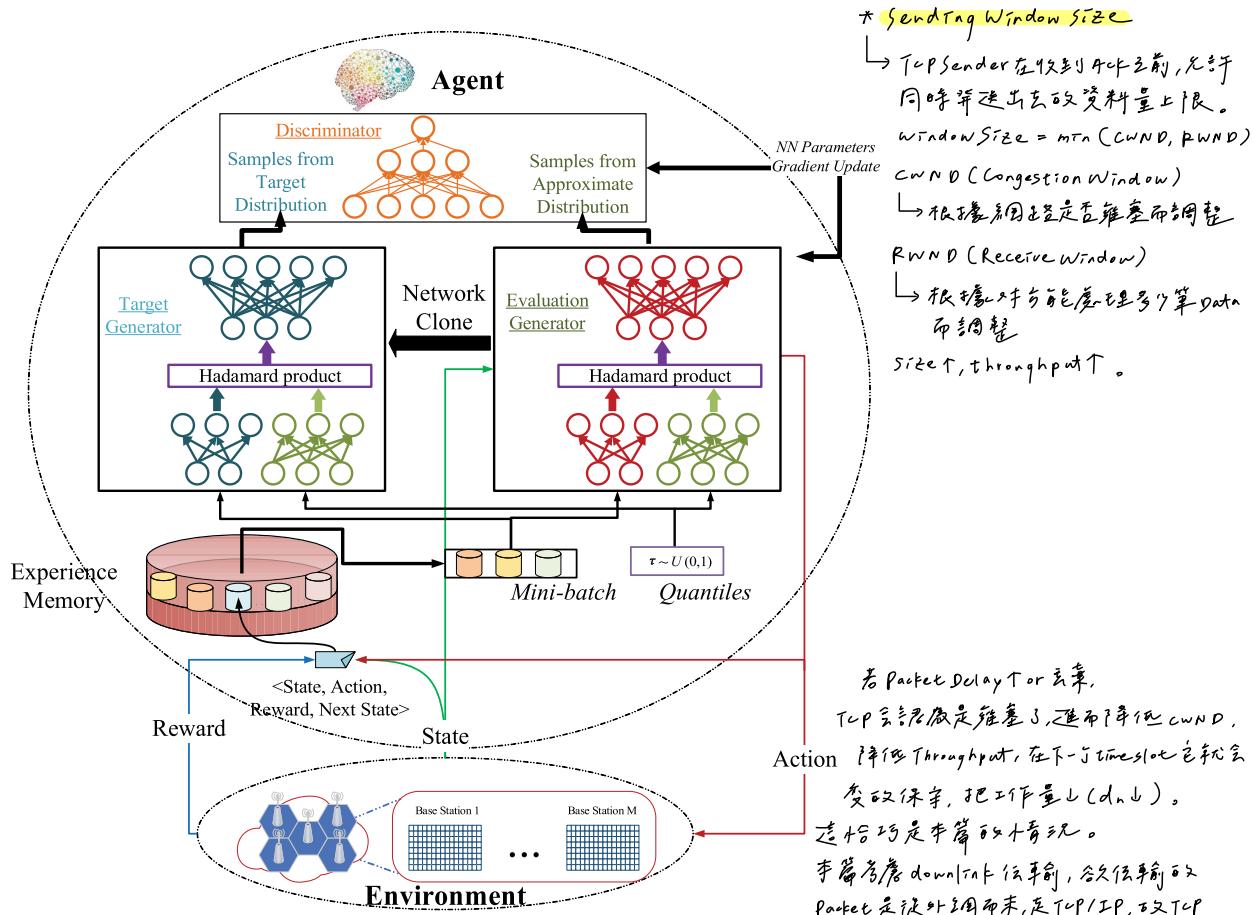


Fig. 2. An illustration of GAN-DDQN for resource management in network slicing.

The bandwidth allocation problem in the RAN network slicing is formulated as follows

$$\begin{aligned} & \max_{w_n} \alpha SE + \sum_{n \in \mathcal{N}} \beta_n \cdot SSR_n \\ & = \max_{w_n} \alpha \left[\frac{\sum_{n \in \mathcal{N}} \sum_{u_n \in \mathcal{U}_n} r_{u_n}}{W} \right] SE \\ & + \sum_{n \in \mathcal{N}} \beta_n \cdot \left[\frac{\sum_{u_n \in \mathcal{U}_n} \sum_{q_{u_n} \in \mathcal{Q}_{u_n}} x_{q_{u_n}}}{\sum_{u_n \in \mathcal{U}_n} |\mathcal{Q}_{u_n}|} \right], \end{aligned} \quad (18)$$

$$\text{s.t. } \sum_{n \in \mathcal{N}} w_n = W, \quad (19)$$

$$\sum_{u_n \in \mathcal{U}_n} |\mathcal{Q}_{u_n}| = d_n, \quad (20)$$

$$x_{p_{u_n}} = \begin{cases} 1, & r_{u_n} \geq \bar{r}_n \text{ and } l_{p_{u_n}} \leq \bar{l}_n, \\ 0, & \text{otherwise.} \end{cases} \quad (21)$$

where α and $\beta = [\beta_1, \beta_2, \dots, \beta_n]$ are the coefficients that adjust the importance of SE and SSR, and β_n refers to the importance weight of SSR_n . In this problem, the objective is to maximize two components: (a) the spectral efficiency (i.e., SE), and (b) the proportion of the packets satisfying the constraint of data rate and latency (i.e., SSR).

Notably, in the current timeslot, d_n depends on both the number of demands and the bandwidth-allocation solution in the previous timeslot, since the (maximum transmission

capacity of RAN belonging to one service is tangled with the provisioning capabilities for this service) For example, (the TCP sending window size is influenced by the estimated channel throughput). Therefore, the traffic demand varies without knowing a prior transition probability, making Eq. (18) difficult to yield a direct solution. However, RL promises to be applicable to tackle this kind of problem. Therefore, we refer to RL to find the optimal policy for network slicing. In particular, consistent with [13], we map the RAN scenario to the context of RL by taking the number of arrived packets in each slice within a specific time window as the state, and the bandwidth allocated to each slice as the action.

III. GAN-POWERED DEEP DISTRIBUTIONAL Q NETWORK

In this section, we describe the proposed GAN-DDQN algorithm, shown in Fig. 2, that address the demand-aware resource allocation problem in network slicing. We then discuss the methods of improving the performance of the algorithm and analyze its convergence.

$$\text{Action Space: } 0 \sim 10 \rightarrow \text{3x3x3x3x3} \rightarrow C^{10+3+3+3+3}$$

$$\text{State: } x+y+z=10, x,y,z \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

A. GAN-DDQN Algorithm

Our previous work [13] has discussed how to apply RL to the resource slicing problem. However, the DQN algorithm used in that work is based on the expectation of the action-value distribution, and thus (does not take into account

Code of Shannon 數字訊號單位推導

$$dB = 10 \cdot \log_{10} \left(\frac{P}{B} \right)$$

$$dBW = 10 \cdot \log_{10} \left(\frac{P}{BW} \right)$$

→ 計算 dBW (相對單位) 轉回 W (絕對單位)

$$A(dBW) = 10 \cdot \log_{10} \left(\frac{P}{BW} \right)$$

$$\hookrightarrow \frac{A}{10} = \log_{10} \left(\frac{P}{BW} \right)$$

$$10^{\frac{A}{10}} = \frac{P}{BW} \Rightarrow P = 10^{\frac{A}{10}} \cdot BW$$

$$BS_tx_power = 10^{dBW}$$

$$rx_power_{dBW} = BS_tx_power_{dBW} - chan_loss_{dB} + UE_rx_gain_{dB}$$

↓ 計算 SNR, 轉回 W

$$rx_power_W = 10^{\left(\frac{rx_power}{10} \right)}$$

rx-power-UE W

$$rate_UE \frac{bits}{s} = band_UE \frac{Hz}{s} * \log_{10} \left(1 + \frac{P_{signal}}{P_{noise}} \right)$$

bits/s

$$(noise-PSD \frac{dBW}{Hz} \xrightarrow{\text{轉為 } W} noise-PSD \frac{W}{Hz}) * band-UE \frac{Hz}{s}$$

W

* Hadamard Product 以外的其他方法

方法	問題/效果
concat	→ 只是把 state 跟 quantile 接起來，沒有做交互作用（平行資訊）
add	→ 必須向量有相同意義與單位，但 state 和 quantile 嵌入意義差很大
<input checked="" type="checkbox"/> Hadamard	→ 遠維交互作用：每一維的特徵都會根據 quantile 調整

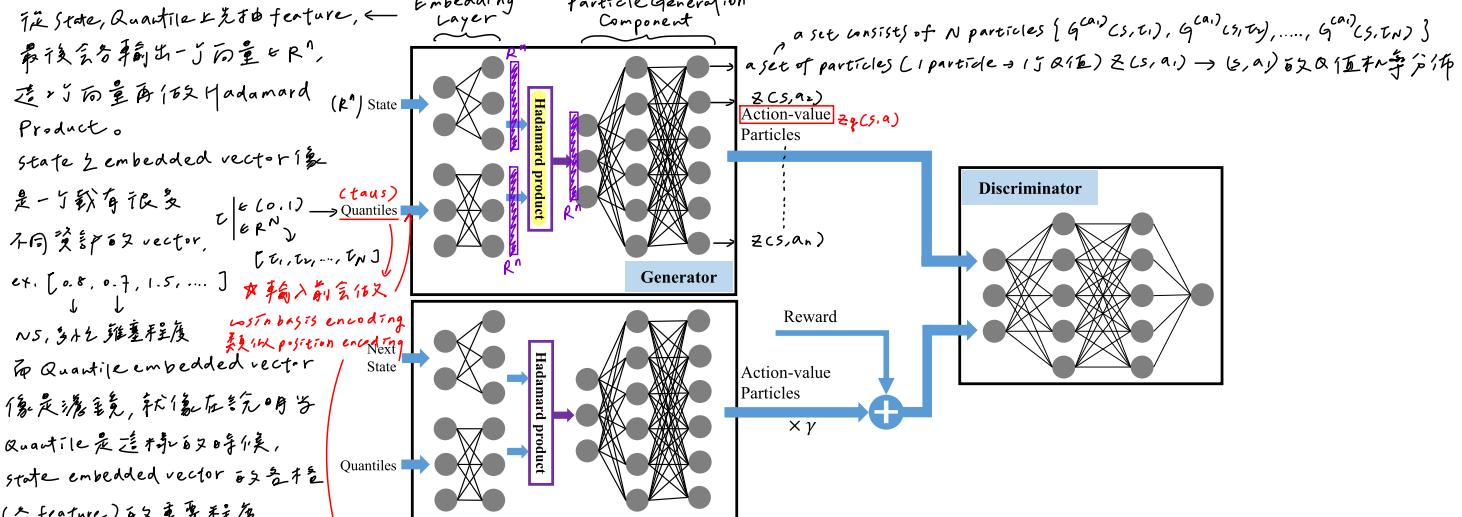
* 以 batch 為度，看 Model 的 input/output shape

Input_shape self.state_size output_shape

⑨ state : (self.batch_size * self.num_samples,)
Quantiles : (self.batch_size * self.num_samples,)
self.embedding_dim

⑩ (self.batch_size, self.num_actions, self.num_samples) (self.batch_size, self.num_actions, self.num_outputs)

$$\text{ex. } a = [a_1, a_2, \dots, a_n], b = [b_1, b_2, \dots, b_n] \Rightarrow a \odot b = [a_1 \cdot b_1, a_2 \cdot b_2, \dots, a_n \cdot b_n]$$

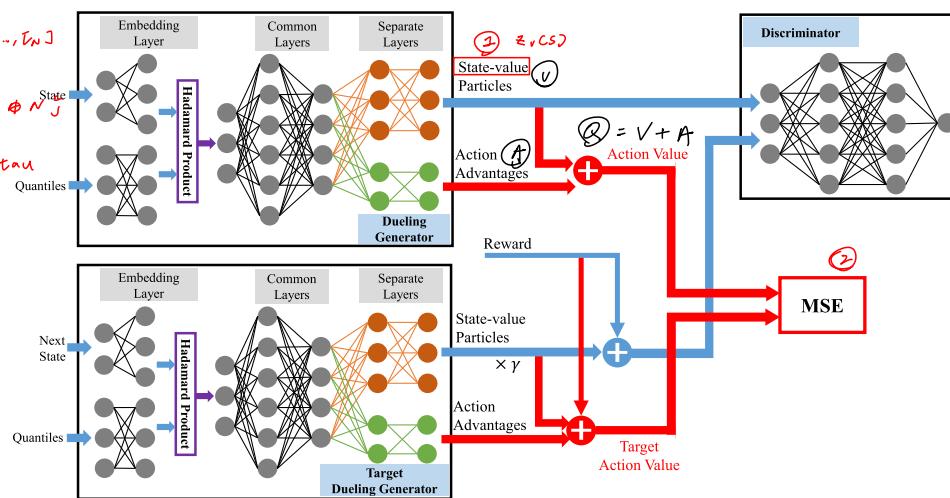


由此, 才会随机地出现两个 R^n , particle & v
这样才会有 Quantile 的粒子理论意义 ($g_{acs,t}$)

一个 state 对应一个 particle。
一个 state 会产生 N 个粒子。
产生 N 个粒子和分布 $\propto N$
不同的 state 会产生不同的粒子

(a) An illustration of GAN-DDQN algorithm.

- 1. Q-value particles \rightarrow State-value particles
- 2. Q-value \oplus TD-Error



(b) An illustration of Dueling GAN-DDQN algorithm.

Fig. 3. The comparison of GAN-DDQN and Dueling GAN-DDQN.

SSR 高发失帧概率 \rightarrow 拼成 9 个分段原图同 DNN 拼成主人图不差细路 - 帧, 路是重叠走训课。

the adverse effects of random noise on the received values of SE and SSR. To overcome this problem, we resort to the combination of the distributional RL and GAN. In this regard, we introduce WGAN-GP to learn the optimal action-value distribution. Specifically, the generator network G outputs a fixed number of samples (we refer to them as *particles* for clarity) that characterize the estimated action-value distribution learned by network G . Similar to [23], we leverage a target generator network \tilde{G} to obtain the target action-value particles. The discriminator network D realizes the 1-Wasserstein criterion when it attempts to minimize the distance between the estimated action-value particles and the target action-value particles calculated by the Bellman optimality operator. GAN-DDQN is able to approximate the optimal action-value distribution by alternately updating networks G and D .

Before we introduce the details of GAN-DDQN algorithm, it is necessary to describe the structure of the networks G and D . Network G consists of three components, which are responsible for state embedding, sample embedding, and particles

generation. The state-embedding and sample-embedding components are both built with two neural layers and process the input state and the quantile samples in parallel. Then, the output of these two components are combined through Hadamard product operation. This step is consistent with [27] to force interaction between the state embedding and sample embedding. Afterwards, the particles generation component, which contains multiple neural layers, takes the fused information as input, outputting several sets of particles where each set is treated as a representation of the corresponding action-value distribution. On the other hand, network D is a multilayer perceptron (MLP) with one neuron in the output layer. Fig. 3(a) further details the structure of GAN-DDQN.

The GAN-DDQN algorithm can be explained as follows, without loss of generality. At iteration t , the agent feeds the current state $S_t = s$ and the samples τ from a uniform distribution (e.g., $U(0,1)$) to network G ; τ is the quantile values of the action-value distribution [27]. Network G outputs a set of estimated action-value particles, denoted as $G(s, \tau)$,

where the particles belonging to action a are denoted as $G^{(a)}(s, \tau)$; the number of $G^{(a)}(s, \tau)$ is N . Then, the agent calculates $Q(s, a) = \frac{1}{N} \sum G^{(a)}(s, \tau), \forall a \in \mathcal{A}$, and selects $a^* = \arg \max_a Q(s, a), \forall a \in \mathcal{A}$ to perform. As a result, the agent receives a reward r , and the environment moves to the next state $S_{t+1} = s'$. The tuple (s, a^*, r, s') is stored into the replay buffer \mathcal{B} . When \mathcal{B} is full, the agent updates networks G and D using all the transition tuples in \mathcal{B} every K iterations.

In the training and updating process, the agent first randomly selects m transitions from \mathcal{B} as a minibatch for training GAN-DDQN. Then, the agent executes the Bellman optimality operator on each transition of the selected minibatch and obtains the target action-value particles. For example, the target action-value particles for the transition i is $y_i = r_i + \gamma \hat{G}^{(a_i^*)}(s'_i, \tau_i)$ where a_i^* is the action with the maximum expectation of action-value particles, i.e., $a_i^* = \arg \max_a \frac{1}{N} \sum \hat{G}^{(a)}(s'_i, \tau_i)$. Finally, the agent uses the following loss functions to train networks D and G , respectively:

$$\begin{aligned} \text{A-B 会走失。} \quad & \mathcal{L}_D = \mathbb{E}_{\substack{\tau \sim U(0,1) \\ (s, a) \sim \mathcal{B}}} [D(G^{(a)}(s, \tau))] \rightarrow A, \text{ 产生假} \\ \text{相当於最大化 B-A,} \quad & - \mathbb{E}_{\substack{(s, a, r, s') \sim \mathcal{B}}}[D(y)] + p(\lambda), \rightarrow B, \text{ 产生真} \\ \text{即目标及真实更明显。} \quad & \text{相当於最大化 E} \quad \leftarrow \mathcal{L}_G = - \mathbb{E}_{\substack{\tau \sim U(0,1) \\ (s, a) \sim \mathcal{B}}} [D(G^{(a)}(s, \tau))] \end{aligned} \quad (23)$$

where $p(\lambda)$ is as mentioned in Eq. (13). The training goal for network D is to increase its accuracy in distinguishing the target action-value particles from the action-value particles produced by network G . The goal of training network G , on the other hand, is to improve its ability to generate the action-value particles that “fool” network D as much as possible. Note that in order to further stabilize the training process, we update the target network \hat{G} every C iterations.

Step by step, we incorporate the aforementioned methods and establish the GAN-DDQN as in Algorithm 1.

B. Convergence Analysis

It has been proven in [24] that the distributional RL can converge when the metric for diverging distributions is p -Wasserstein distance. On the other hand, the fundamental guidance for distinguishing the target and estimated distributions in WGAN-GP is 1-Wasserstein distance. Therefore, the convergence of GAN-DDQN can be analyzed from the perspective of WGAN-GP’s convergence on the data sampled from the dynamic RL interaction process. However, as explored in [46], in many currently popular GAN architectures, converging to the target distribution is not guaranteed and oscillatory behavior can be observed. This is a twofold challenge for GAN-DDQN, as we must ensure both the stationarity of the target distribution and the convergence of the WGAN-GP to this target distribution.

In an idealized WGAN-GP, the generator should be able to learn from the target distribution, and the discriminator should be able to learn any 1-Lipschitz function to produce the exact Wasserstein distance. However, the target distribution

☆ 真实连接分布 Diving GAN-DDQN - 不容易训练 5 次才训练出来

☆ G, 为了避免 mode collapse

WGan-GP + Distributional DQN ☆ er 1e-4

Algorithm 1 GAN-DDQN

- 1: Initialize a generator G and a discriminator D with random weights θ_G and θ_D respectively, the number of particles N , gradient penalty coefficient λ , batch size m , discount factor γ .
- 2: Initialize a target generator \hat{G} with weight $\theta_{\hat{G}} \leftarrow \theta_G$, a replay buffer $\mathcal{B} \leftarrow \emptyset$, the iteration index $t = 0$.
- 3: **repeat**
- 4: The agent observes $S_t = s$.
- 5: The agent samples $\tau \sim U(0, 1)$.
- 6: The agent calculates $Q(s, a) = \frac{1}{N} \sum G^{(a)}(s, \tau), \forall a \in \mathcal{A}$.
- 7: The agent performs $a^* \leftarrow \arg \max_a Q(s, a), \forall a \in \mathcal{A}$.
- 8: The agent receives the system utility J and observes $S_{t+1} = s'$.
- 9: The agent performs the reward-clipping with respect to J and gets the reward r .
- 10: The agent stores transition (s, a^*, r, s') in \mathcal{B} .
- 11: If \mathcal{B} is full, the agent updates the weights of network G and network D every K iterations. $\rightarrow K=2000$ (Learning windows)
- 12: **# Train GAN**
- 13: **repeat**
- 14: The agent samples a minibatch $\{s, a, r, s'\}_{i=1}^m$ from \mathcal{B} without replacement. \rightarrow minibatch 中无任何重複之 sample
- 15: The agent samples a minibatch $\{\tau\}_{i=1}^m \sim U(0, 1)$.
- 16: The agent gets the target action-value particles $y_i = r_i + \gamma \hat{G}^{(a_i^*)}(s'_i, \tau_i)$, where the optimal action is $a_i^* = \arg \max_a \frac{1}{N} \sum \hat{G}^{(a)}(s'_i, \tau_i), \forall a \in \mathcal{A}$.
- 17: The agent samples a minibatch $\{\varepsilon\}_{i=1}^m \sim U(0, 1)$, and sets $\hat{x}_i = \varepsilon_i y_i + (1 - \varepsilon_i) G^{(a_i)}(s_i, \tau_i)$. \rightarrow Gradient Penalty \rightarrow $\frac{1}{m} \sum_{i=1}^m \mathcal{L}_i$, where $\mathcal{L}_i = D(G^{(a_i)}(s_i, \tau_i)) - D(y_i) + \lambda (\|\nabla_{\hat{x}_i} D(\hat{x}_i)\|_2 - 1)^2$.
- 18: The agent updates the weights θ_D by leveraging gradient descent algorithm to $\frac{1}{m} \sum_{i=1}^m \mathcal{L}_i$, where $\mathcal{L}_i = D(G^{(a_i)}(s_i, \tau_i)) - D(y_i) + \lambda (\|\nabla_{\hat{x}_i} D(\hat{x}_i)\|_2 - 1)^2$.
- 19: The agent updates the weights θ_G by leveraging gradient descent algorithm to $-\frac{1}{m} \sum_{i=1}^m D(G^{(a_i)}(s_i, \tau_i))$.
- 20: **until** All the transitions in \mathcal{B} are used for training.
- 21: The agent clones network G to the target network \hat{G} every C iterations by resetting $\theta_{\hat{G}} = \theta_G$.
- 22: The iteration index is updated by $t \leftarrow t + 1$.
- 23: **until** A predefined stopping condition (e.g., the $\frac{1}{m} \sum_{i=1}^m \mathcal{L}_i$, the preset number of iterations, etc.) is satisfied.

will not be stationary as the target network \hat{G} regularly updates its weights; thus an idealized WGAN-GP might not be successful in practice. Fortunately, a slight change in the target distribution has little effect on the convergence of WGAN-GP. For example, suppose the real distribution that is the ultimate learning goal of WGAN-GP is a Gaussian distribution with a mean of 100 and a standard deviation of 1, and suppose that the target distribution that the WGAN-GP is expected to approximate at each stage is a Gaussian distribution with a standard deviation of 1 and a mean that starts at 0, increasing periodically by $\Delta\mu$. WGAN-GP will need more updates to learn the target distribution if $\Delta\mu$ is large; the number of

△ 从 Target distribution 3 到 5 stage 变化平缓。

updates is difficult to determine. However, if $\Delta\mu$ is small, a few times of updates is sufficient for WGAN-GP to learn the changed target distribution. Hence, the small $\Delta\mu$ is more potential to enable WGAN-GP to learn the real distribution smoothly. ↗ 附录 Dirac-WGAN-GP 的 Theorem 2, 由此推出

To analyze the convergence characteristic of WGAN-GP while avoiding directly dealing with sophisticated data and WGAN-GP model, [46] introduces a simple but illustrative model, namely **Dirac-WGAN-GP**. Specifically, Dirac-WGAN-GP consists of a linear discriminator $D_\psi(x) = \psi \cdot x$ and a generator with parameter θ that indicates the position of the Dirac distribution (i.e., δ_θ) output by the generator. Whilst the real data distribution p_d is given by a Dirac-distribution concentrated at ξ (i.e., δ_ξ). It is worthy to further investigate the training characteristic of Dirac-WGAN-GP when the real data distribution (i.e., δ_ξ) is varying during the training process, like the typical situation in RL. Consistent with [46], we carry out analysis based on Dirac-WGAN-GP, and we have the following Theorem 1

Theorem 1: When trained with gradient descent with a fixed number of the generator and the discriminator updates and a fixed learning rate $h > 0$, if the value of ξ varies dramatically, Dirac-WGAN-GP needs more learning steps to converge from the old optimal boundary to the new one after the variation of ξ .

We leave the proof of Theorem 1 in Appendix. As for WGAN-GP, we further have the following Corollary 1

Corollary 1: WGAN-GP could converge to the optimal boundary more rapidly if the real data change by a small amount.

Corollary 1 reveals that estimating the optimal action-value distribution requires a large amount of training if we directly use the system utility as the reward in RL. Therefore, we put forward a new **reward-clipping** mechanism to prevent the target action-value distribution from greatly changing. Specifically, (assuming that there are T thresholds that partition the system utility, we set $T+1$ constants whose values are much smaller than the system utility) Then the system utility can be clipped to these $T+1$ constants that are taken as the rewards in RL. For example, if $T=2$ and the clipping constants are $-\eta$, 0 , and η ($\eta > 0$), then the clipping strategy can be formulated by Eq. (24), where c_1 and c_2 ($c_1 > c_2$) are the manually set thresholds: ↗ if $J(w, d) \geq c_1$, then $r = \eta$; if $J(w, d) \leq c_2$, then $r = -\eta$; if $c_2 < J(w, d) < c_1$, then $r = 0$.

$$r = \begin{cases} \eta, & J(w, d) \geq c_1, \\ 0, & c_2 < J(w, d) < c_1, \\ -\eta, & J(w, d) \leq c_2. \end{cases} \quad (24)$$

However, as T becomes larger, the number of the manually set parameters in the reward-clipping mechanism increases, which makes the parameter setting process more sophisticated. Therefore, we adopt the reward-clipping mechanism defined in Eq. (24) as an experiment. Note that introducing the reward-clipping mechanism to GAN-DDQN algorithm is easy, and we only need to apply the reward-clipping mechanism to the system utility before storing the transition tuple in the replay buffer, which is described in line 9 of Algorithm 1. ↗ 平衡的双端 Q & A (T=1)

↗ if $J(w, d) \geq c$ reward = 1
else reward = 0

→ 双端 WGAN-GP 的 Corollary 2: **在 Action Space 1 $\times N_1$ particles, 但只用其中 N particles 训练**

The training of **GAN-DDQN** is not a trivial task since it uses the data yielded from a dynamic environment, and only a tiny portion of the output of the generator is useful for gradient calculation. One intuitive indicating to alleviate the training problem is to carefully adjust the values of GAN-DDQN's hyper-parameters, such as the discount factor γ , the gradient penalty coefficient λ , etc. Nevertheless, we plan to make systemic and architectural changes to the generator and the loss function. Particularly, inspired by [36], which uses a **specialized dueling Q network** to separate the action value into a state-value stream and an advantage stream, we divide the approximation of the action-value distribution into the approximation of the state-value distribution and the approximation of the advantage function for each action. This dueling architecture ignores the trivial variations of the environment and focuses on some crucial states to enhance the stability of DRL algorithms [35]. In addition, in our improved model, namely **Dueling GAN-DDQN**, the **loss function** of the **discriminator** turns to work on the estimated and target state-value distributions. Moreover, the squared TD error is added to the **generator's loss** as the criterion that measures the distance of the estimated and target action values. ↗ GAN-DDQN 与 G^2C�, 这是 G^2C�

The detailed structure of Dueling GAN-DDQN is presented in Fig. 3(b), and we remarkably highlight the key differences from GAN-DDQN. It can be observed that the significant difference compared with GAN-DDQN is the generator or the dueling generator for the sake of distinguishing. In the dueling generator, after Hadamard product operation, we continue to handle the output using multiple neural layers (the **common layers**). Then, the refined information is separated into two paths, one flowing to a neural network to approximate the state-value distribution, and the other flowing to another neural network to estimate the action advantage function. Accordingly, the dueling generator outputs not only particles from the approximated state-value distribution but also the estimated action advantage values for each action. Note that the discriminator of Dueling GAN-DDQN has the same structure as GAN-DDQN.

Similarly to our analysis of the random variable Z_v^π , we analyze the random variable $Z_v^\pi(s)$, which denotes the return obtained by following a policy π from state s . Then we have

$$V^\pi(s) = \mathbb{E}[Z_v^\pi(s)], \quad (25)$$

and an analogous distributional Bellman equation for Z_v

$$Z_v^\pi(s) \stackrel{D}{=} \mathbb{E}_{\substack{a \sim \mathcal{A} \\ s' \sim \mathcal{S}}} [R + \gamma Z_v^\pi(s')]. \quad (26)$$

It is difficult to find the distributional Bellman optimality operator for Z_v^π . Even worse, Eq. (26) indicates that the iterative calculation of Z_v^π requires a reward from every state-action pair, which is a noticeable time-consuming operation. Therefore, we introduce a degraded but simplified method to estimate Z_v^π , which is to minimize the difference between the estimated Z_v^π and $\mathcal{T}Z_v^\pi$ calculated by

$$\mathcal{T}Z_v^\pi \stackrel{D}{=} r + \gamma Z_v^\pi(s'), \quad (27)$$

where r and s' are from the transition (s, a, r, s') sampled from the replay buffer. This degraded approximation may fail to find the optimal state-value distribution, yet it can significantly reduce computation time. In addition, only considering the 1-Wasserstein loss for the state-value distribution results in the network G weights related to the action advantage function not being trained. Therefore, we leverage the TD error to measure the difference of the estimated and the target action values, where the action value is calculated by adding the corresponding action advantage value to the mean of the state-value particles. As a consequence, the ultimate loss function for training **Dueling GAN-DDQN** is composed of the 1-Wasserstein distance and the squared TD error, which can be formulated as follows

$$\begin{aligned} \mathcal{L}_D = & \mathbb{E}_{\substack{\tau \sim U(0,1) \\ (s, a) \sim \mathcal{B}}} [D(G_v(s, \tau))] \\ & - \mathbb{E}_{\substack{\tau \sim U(0,1) \\ (r, s') \sim \mathcal{B}}} [D(r + \gamma \hat{G}_v(s', \tau))] + p(\lambda), \end{aligned} \quad (28)$$

$$\mathcal{L}_G = - \mathbb{E}_{\substack{\tau \sim U(0,1) \\ (s, a) \sim \mathcal{B}}} [D(G_v(s, \tau))] + \frac{1}{2} \zeta^2, \quad (29)$$

where G_v denotes the state-value particles output by the dueling generator, and ζ^2 is the squared TD error as defined in Eq (7). Algorithm 2 and Fig. 3(b) provide the details of Dueling GAN-DDQN.

IV. SIMULATION RESULTS AND NUMERICAL ANALYSIS

A. Simulation Environment Settings

In this part, we verify the performance of **GAN-DDQN** and **Dueling GAN-DDQN** in a RAN scenario where there are three types of services (i.e., VoLTE, video, and URLLC) and three corresponding slices in one serving BS, as in [13]. There exist 100 registered subscribers randomly located within a 40-meter radius circle surrounding the BS. These subscribers generate standard service traffics as summarized in Table II based on 3GPP TR 36.814 [47] and TS 22.261 [48]. The total bandwidth is 10 MHz, and the bandwidth allocation resolution is 1 MHz or 200 KHz. We will show the simulation results for both cases. On the other hand, the packet size of URLLC service has a strong influence on the system utility. For example, it is difficult to meet the latency requirement of URLLC service when the packet size is large, if there is insufficient bandwidth guaranteed for transmission. As a result, SSR degrades, and the system utility is reduced. Therefore, we simulate the network slicing scenario with suitably-sized URLLC packets.

With the mapping shown in Table III, RL algorithms can be used to optimize the system utility (i.e., the weighted sum of SE and SSR). Specifically, we perform round-robin scheduling within each slice at 0.5 ms granularity; that is, we sequentially allocate the bandwidth of each slice to the active users within each slice every 0.5 ms. Besides, we adjust the bandwidth allocation to each slice per second. Therefore, the agent updates its neural networks every second. Considering the update interval of the bandwidth-allocation process is much larger than the service arrival interval, the number of arrived

* lr % 级别 1e-3

Algorithm 2 Dueling GAN-DDQN

- 1: Initialize a dueling generator G and a discriminator D with random weights θ_G and θ_D respectively, the number of particles N , gradient penalty coefficient λ , batch size m , discount factor γ , $n_{critic} = 5$.
- 2: Initialize a target dueling generator \hat{G} with weight $\theta_{\hat{G}} \leftarrow \theta_G$, a replay buffer $\mathcal{B} \leftarrow \emptyset$, the iteration index $t = 0$.
- 3: **repeat**
- 4: The agent observes $S_t = s$.
- 5: The agent samples $\tau \sim U(0, 1)$.
- 6: The agent feeds s and τ to network G , getting the state-value particles $G_v(s, \tau)$ and each action advantage value $G_{ad}^{(a)}(s, \tau)$, $\forall a \in \mathcal{A}$. mean of vs
- 7: The agent calculates $V(s) = \frac{1}{N} \sum G_v(s, \tau)$.
- 8: The agent calculates $Q(s, a) = V(s) + G_{ad}^{(a)}(s, \tau)$, $\forall a \in \mathcal{A}$.
- 9: The agent performs $a^* \leftarrow \arg \max_a Q(s_t, a)$.
- 10: The agent receives the system utility J and observes a new state $S_{t+1} = s'$.
- 11: The agent performs the reward-clipping with respect to J and gets the reward r .
- 12: The agent stores transition (s, a^*, r, s') in \mathcal{B} .
- 13: **# Train network D**
- 14: **for** $n = 1$ to n_{critic} **do**
- 15: The agent randomly samples $\{s, a, r, s'\}_{i=1}^m$ from \mathcal{B} .
- 16: The agent samples $\{\tau\}_{i=1}^m$ and $\{\varepsilon\}_{i=1}^m$ from $U(0, 1)$.
- 17: The agent gets $y_i = G_v(s_i, \tau_i)$, and $\hat{y}_i = r_i + \gamma \hat{G}_v(s'_i, \tau_i)$.
- 18: The agent sets $\hat{x}_i = \varepsilon_i \hat{y}_i + (1 - \varepsilon_i) y_i$.
- 19: The agent updates the weights θ_D by leveraging gradient descent algorithm to $\frac{1}{m} \sum_{i=1}^m \mathcal{L}_i$, where $\mathcal{L}_i = D(y_i) - D(\hat{y}_i) + \lambda (\|\nabla_{\hat{x}_i} D(\hat{x}_i)\|_2 - 1)^2$.
- 20: **end for**
- 21: **# Train network G**
- 22: The agent randomly samples $\{s, a, r, s'\}_{i=1}^m$ from \mathcal{B} .
- 23: The agent samples $\{\tau\}_{i=1}^m$ from $U(0, 1)$.
- 24: The agent calculates the estimated action value $Q_i = \frac{1}{N} \sum G_v(s_i, \tau) + G_{ad}^{(a_i)}(s_i, \tau)$, and the target action value $\hat{Q}_i = r_i + \gamma \frac{1}{N} \sum G_v(s'_i, \tau) + \gamma \max_a G_{ad}^{(a)}(s'_i, \tau)$, $\forall a \in \mathcal{A}$.
- 25: The agent updates the weights θ_G by leveraging gradient descent algorithm to $\frac{1}{m} \sum_{i=1}^m [-D(G_v(s_i)) + \frac{1}{2} (\hat{Q}_i - Q_i)^2]$.
- 26: The agent clones network G to the target network \hat{G} every C iterations by resetting $\theta_{\hat{G}} = \theta_G$.
- 27: The iteration index is updated by $t \leftarrow t + 1$.
- 28: **until** A predefined stopping condition (e.g., the $\frac{1}{m} \sum_{i=1}^m \mathcal{L}_i$, the preset number of iterations, etc.) is satisfied.

→ packet 速率 及时序 比上层 低很多 且很慢。
→ 不会有一个 NS 中的 packet 速率过快导致拥塞。且这个带宽限制
packets (i.e., the state) is rarely zero when updating the agents.
Therefore, it is reasonable to ignore the situation of zero bandwidth for any NS. Moreover, this filter setting narrows the range for the action exploration, as well as enhancing the

eMBB 採 Pareto 分佈 ($\text{base} = 1.2$)

Pareto 分布是一種 **heavy-tailed** 分布，用來模擬：

- 大量小事件、少數大事件的情況
- 比如：大多數人很少上網，但有一小部分人流量爆炸高（比如一直在看 4K 影片）

Pareto(α) 的特性是：

- 越小的 α （例如 $\alpha = 1.2$ ），尾巴越粗（也就是高機率出現極大值）
- 它適合用來模擬 video streaming、web traffic 等 bursty 行為

為什麼乘上 6ms？

python

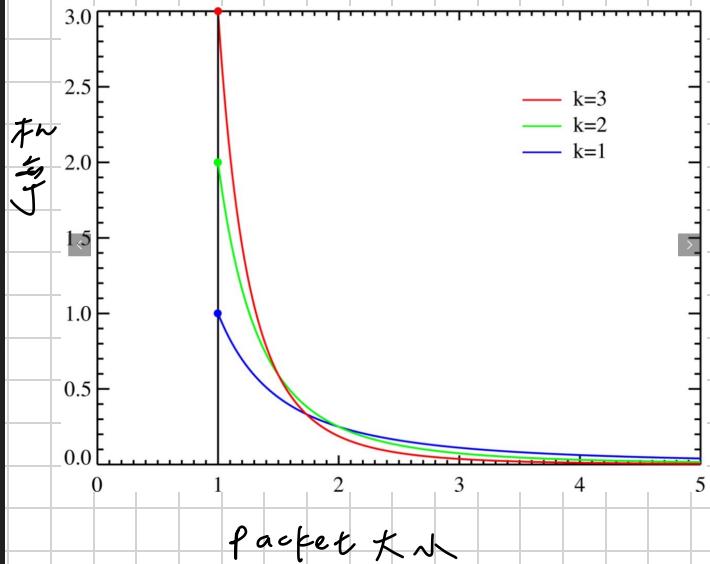
Copy Edit

```
tmp_readtime = pareto(1.2) * 6 ms
```

代表：一個 burst 平均持續時間大約是幾毫秒（這裡以 6ms 為單位），但偶爾會出現極大的時間間隔。

總結這行的物理意義：

為每個 eMBB UE 建立一個隨機「封包產生時間間隔」，模擬影片串流中偶爾產生的 bursty traffic（連續出現資料）、或 silent interval（較長時間沒封包）。



URLLC 採 Exponential 分佈

✓ 指數分布的特性：

- 記憶無關 (memoryless)：未來的事件發生機率與過去無關
- 常用來模擬隨機突發事件發生的間隔
- PDF 為： $f(x) = \lambda e^{-\lambda x}$ ，平均為 $1/\lambda$

為什麼適合 URLLC？

URLLC (Ultra-Reliable Low-Latency Communications) 通常是：

- 突發的事件觸發（例如：感測器發送緊急資訊、車輛碰撞偵測）
- 並不是週期性的傳輸
- 所以封包間隔時間應該是隨機且稀疏的

使用指數分布可以很自然地模擬：

URLLC 封包是在一段平均為 180ms 的時間內突發發生。

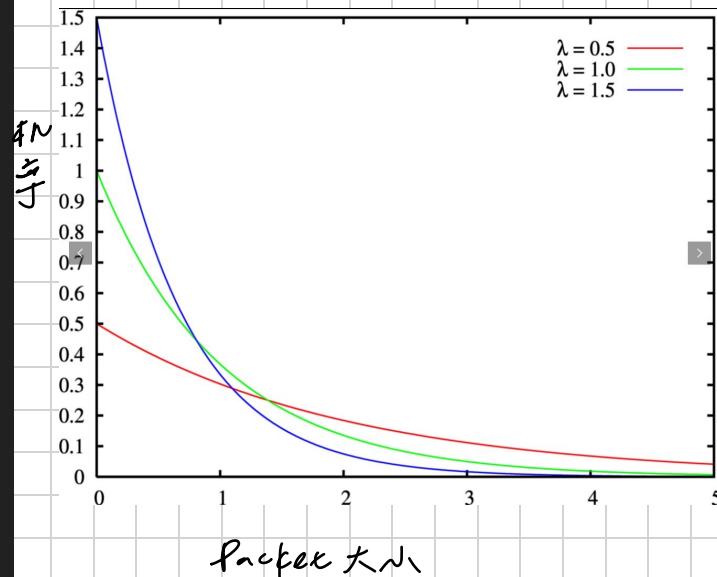


TABLE II
A BRIEF SUMMARY OF KEY SETTINGS FOR TRAFFIC GENERATION PER SLICE

	语音电话	视频 (eMBB)	工业控制
Bandwidth	20 MHz	Video (eMBB)	URLLC
Scheduling	Round robin per slot (0.5 ms)		
Slice Band Adjustment (Q-Value Update)	1 second (2000 scheduling slots)		
Channel	Rayleigh fading → only consider large-scale fading		
User No. (100 in all)	46	46	8
Distribution of Inter-Arrival Time per User	Uniform [Min = 0, Max = 160ms]	Truncated stationary distribution [Exponential Para = 1.2, Mean = 6 ms, Max = 12.5 ms]	Exponential [Mean = 180 ms]
Distribution of Packet Size	Constant (40 Byte)	Truncated Pareto [Exponential Para = 1.2, Mean = 100 Byte, Max = 250 Byte]	Variable constant: {6.4, 12.8, 19.2, 25.6, 32} KByte or {0.3, 0.4, 0.5, 0.6, 0.7} MByte
SLA: Rate	51 Kbps	100 Mbps	10 Mbps
SLA: Latency	10 ms	10 ms	1 ms

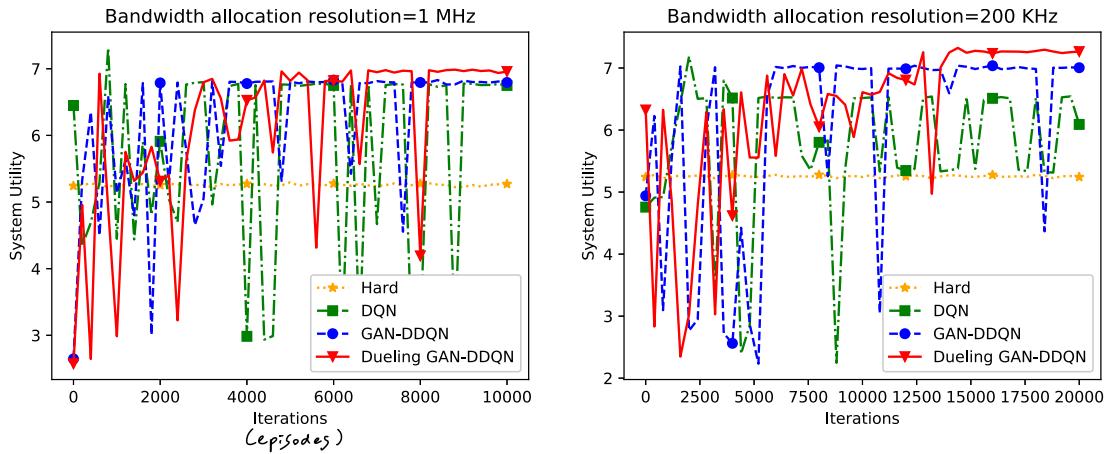


Fig. 4. An illustration of performance comparison between different slicing schemes (hard slicing, DQN, GAN-DDQN, Dueling GAN-DDQN).

TABLE III
THE MAPPING FROM RESOURCE MANAGEMENT FOR NETWORK SLICING TO RL ENVIRONMENT

RL Environment	Radio Resource Slicing
State	The number of arrived packets in each slice within a specific time window
Action	bandwidth allocation to each slice
Reward	Clipped weighted sum of SE and SSR in 3 sliced bands

stability of the training process. Meanwhile, this filter setting does affect our main results.

B. Simulation Results

In this part, we show the simulation results of the proposed **GAN-DDQN** and **Dueling GAN-DDQN** algorithms, in comparison with the **hard slicing method** and the **standard DQN-based** scheme. Hard slicing means that each service is always allocated with $\frac{1}{3}$ of the whole bandwidth (because there are three types of services in total); round-robin scheduling is conducted within each slice. The DQN-based bandwidth allocation scheme was first proposed in [13], which directly

applied the original DQN algorithm [23] to the network slicing scenario. Notably, our previous works in [13] have demonstrated the classical DQN-driven method is superior to other machine learning methods (e.g., long short-term memory (LSTM)-based prediction-before-allocation method). Therefore, due to the space limitation, we put more emphasis on the performance comparison with the classical DQN in [13].

①) **Small Packets for URLLC Service:** We first observe the performance of the proposed algorithms within the scenario that the packet size of URLLC service is small. The traffic parameters are shown in Table II. We consider two cases: the **bandwidth allocation resolution** is either **1 MHz** or **200 KHz**. The importance weights in the optimization objective (i.e., Eq (18)) are set to $\alpha = 0.01, \beta = [1, 1, 1]$. The values of the clipping parameters c_1 and c_2 are determined heuristically.⁴ In both cases, we set $c_1 = 6.5, c_2 = 4.5$ to clip the system utility according to Eq. (24), where η is fixed at 1. The experimental evaluation of the reward-clipping setting is investigated hereinafter. Fig. 4 depicts the variations of the

⁴We first directly regard the system utility as the reward in order to find the range of the system utility, and then try different combinations of the two parameters (i.e., c_1 and c_2) to find the suitable values that guarantee both performance and stability.

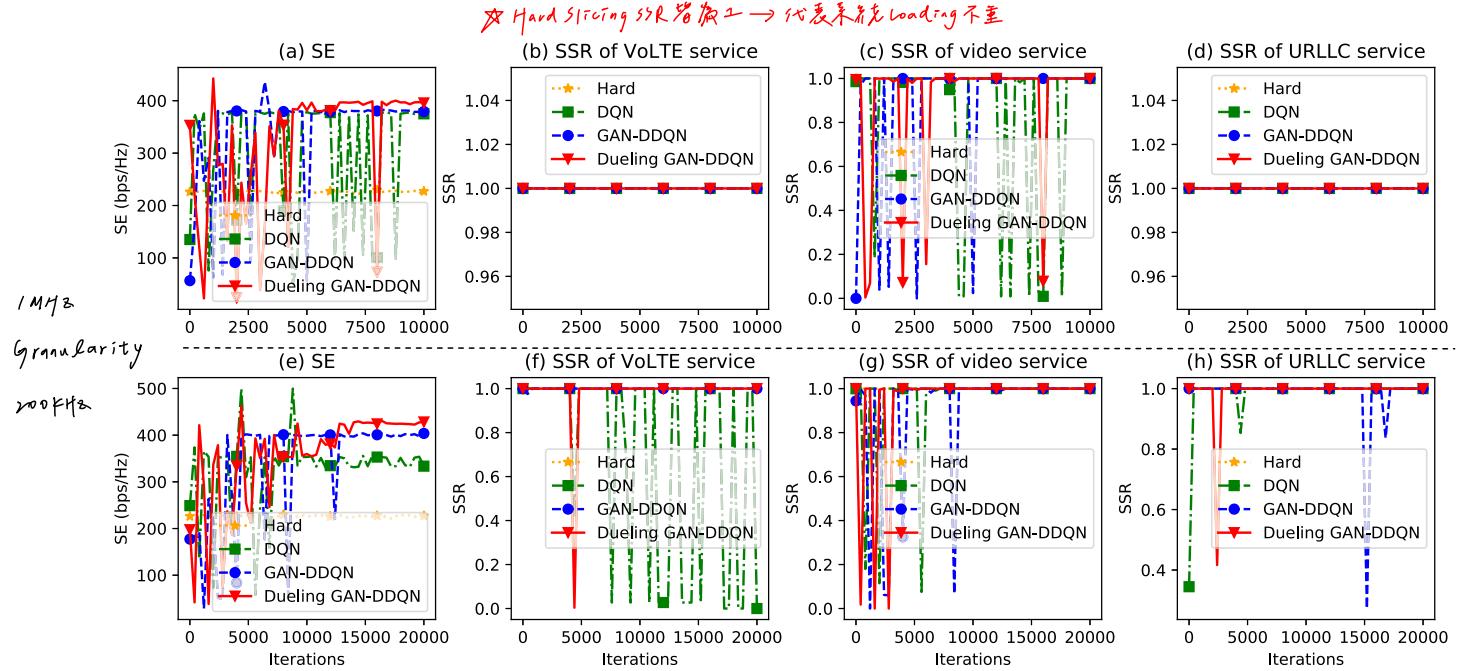


Fig. 5. An illustration of SE and SSR in the different cases where the bandwidth allocation resolution is 1 MHz (shown in the top sub-figures) and 200 KHz (shown in the bottom sub-figures).

system utility with respect to the iteration index. The left part of Fig. 4 shows that when the bandwidth allocation resolution is 1 MHz, the three RL-based algorithms perform similarly, but Dueling GAN-DDQN is slightly better; DQN is the most erratic in training. The right part of Fig. 4 illustrates that when the bandwidth allocation resolution becomes 200 KHz, GAN-DDQN and Dueling GAN-DDQN expand the gap to DQN, which demonstrates the performance improvement coming from distributional RL by the characterization of the action-value or state-value distributions. It's worth noting that Dueling GAN-DDQN improves visibly over GAN-DDQN in both performance and stability, consistent with our previous discussion. Furthermore, it can be observed in Fig. 4 that the system utility obtained by the three RL-based algorithms is significantly greater than that for the hard slicing scheme. The reason for this lies in that the RL-based algorithms can dynamically and reasonably manage bandwidth resources, thereby avoiding wasted resources and improving resources utilization. Moreover, the utilization of the bandwidth resource further gets improved when we slice the bandwidth more finely. Fig. 4 shows that the system utility obtained by the two GAN-based algorithms, especially Dueling GAN-DDQN, becomes significantly larger when the bandwidth allocation resolution changes to 200 KHz from 1 MHz, but that the performance of DQN is slightly degraded.

Fig. 5 presents the variations of SE and SSR with respect to the iteration index for both bandwidth allocation resolution settings (i.e., 1 MHz and 200 KHz). It can be observed from Fig. 5 that SE curves are basically consistent with the system utility curves. However, the SSR curves of the three algorithms for the NSs show different patterns. When the bandwidth allocation resolution is 1 MHz, the SSRs for both VoLTE and URLLC services reach 100% with iterative training.

$$J = \alpha \cdot SE + \sum_{n \in N} \beta_n \cdot SSR_n$$

$$\alpha = 0.01, \beta = [1, \dots, 1]$$

Nevertheless, for the SSR of video service, GAN-DDQN and Dueling GAN-DDQN basically converge to 100% after 5000 iterations, while DQN shows no obvious signs of convergence. When the bandwidth allocation resolution is 200 KHz, it can be observed that GAN-DDQN and Dueling GAN-DDQN, by and large, realize 100% of SSR for all three services by the end of training, but DQN shows extreme instability for VoLTE service. Note that the unusual sudden performance drop late in training is caused by the tiny nonzero exploration rate in the ϵ -greedy exploration strategy.

In Fig. 6, we illustrate the policy learned by the three algorithms when the bandwidth allocation resolution is 200 KHz. It can be observed that all three algorithms converge after 15000 iterations. However, there are some differences between the learned bandwidth allocation policies. The DQN agent allocates the least bandwidth to VoLTE service and keeps it unchanged; as a result, SSR of VoLTE service does not always reach 100%. Between the GAN-DDQN agent and the Dueling GAN-DDQN agent, the latter behaves more intelligently, which is prominently manifested in the fact that it maximizes the bandwidth allocated to video service and reduces the bandwidth allocated to the other two services while meeting the SLA. The Dueling GAN-DDQN agent provides as much bandwidth as possible to satisfy the SLA of video service which is requested frequently and bandwidth-consuming, thus improving the SE. Besides, Dueling GAN-DDQN agent provides a policy to better balance the demands of VoLTE and URLLC services that are relatively rarely requested.

We next investigate the impact of the reward-clipping mechanism. Fig. 7 shows the differences in system utility during the iterative learning of GAN-DDQN with and without the reward clipping when the bandwidth allocation resolution is 1 MHz. When there is no reward clipping, GAN-DDQN directly takes

(4-1)

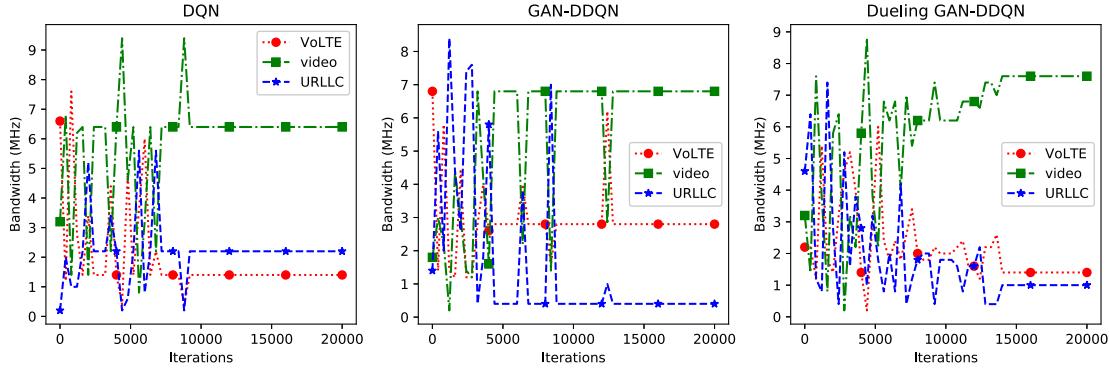


Fig. 6. An illustration of bandwidth allocation schemes in the case where the bandwidth allocation resolution is 200 KHz.

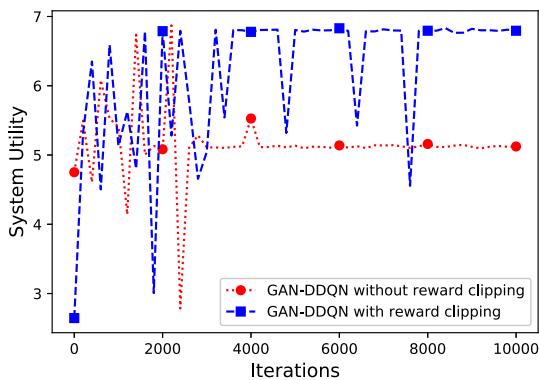


Fig. 7. Comparison of GAN-DDQN with and without reward clipping.

the system utility as the reward. Note that the values of system utility, although they fluctuate, are much larger than the clipping constants set manually in the reward-clipping mechanism. As a result, if the system utility is directly used as the reward, the target action-value distribution might vary significantly, making it difficult for GAN-DDQN to converge. Therefore, GAN-DDQN without reward clipping requires more training steps to converge from one equilibrium to a new one. It can be observed from Fig. 7 that the GAN-DDQN performs significantly better with reward clipping than without. The simulation results verify the effectiveness of GAN-DDQN together with reward clipping.

2) *Large Packets for URLLC Service*: In this part, we consider the case where the packet size of URLLC service is evenly sampled from $\{0.3, 0.4, 0.5, 0.6, 0.7\}$ MByte, which gives a considerably larger packet size than we just analyzed and requires more bandwidth to guarantee meeting the SLA of URLLC service. In the case where bandwidth allocation resolution is 1 MHz, we set $c_1 = 5.7$, $c_2 = 3$ to clip the system utility according to Eq. (24), where η is fixed to 1. Fig. 8 shows the performance of each slice algorithm, from which it can be observed that **Dueling GAN-DDQN** is way ahead of the others in terms of system utility. However, quite unexpectedly, **DQN** performs poorly, worse even than hard slicing scheme. Fig. 9 reveals the details of SE and SSR, from which we can find that **Dueling GAN-DDQN** agent learned a policy that maximizes system utility by sacrificing the SSR of URLLC service in

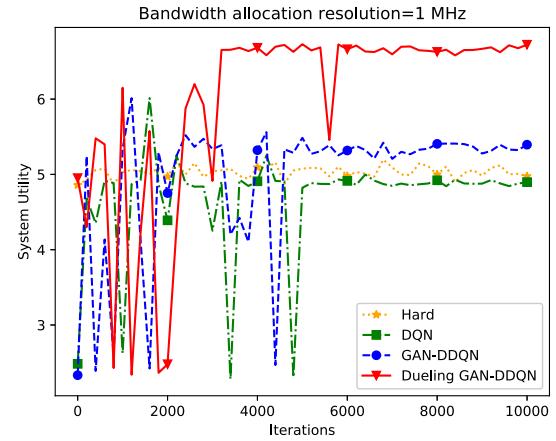


Fig. 8. An illustration of performance comparison between different slicing schemes in the case where the packets of URLLC service are large and the bandwidth allocation resolution is 1 MHz.

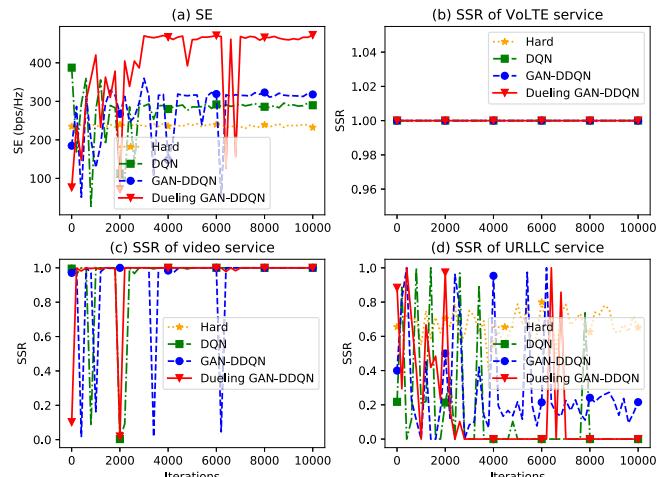


Fig. 9. An illustration of SE and SSR in the case where the bandwidth allocation resolution is 1 MHz.

exchange for higher SE. The reason for this is that when SSR is equally important to all three services (i.e., $\beta = [1, 1, 1]$), it is challenging for URLLC service to satisfy its SLA given the large transmission volume and the strictly low latency requirement. Therefore, we further investigate the situation in

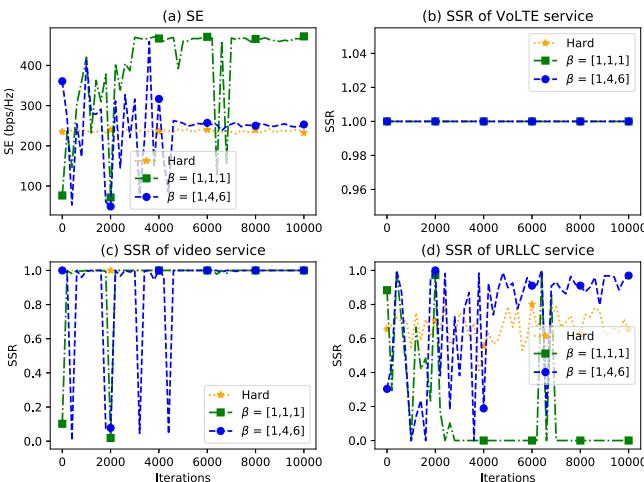


Fig. 10. An illustration of SE and SSR achieved by **Dueling GAN-DDQN** when the bandwidth allocation resolution is 1 MHz and the importance weights of SSR are [1, 1, 1] and [1, 4, 6], respectively.

which URLLC service is more concerned while keeping video service dominating VoLTE service, which is reflected in the value of β changing to [1, 4, 6] from [1, 1, 1]. Fig. 10 presents the results for these different β settings and demonstrates that the adjusted importance weight makes the SLA of URLLC service well guaranteed; meanwhile, the SLA of the other two services can be 100% satisfied as well. However, the requests for URLLC service are scarce in a timeslot—in our simulation, it is the interval between two updates of the agent, defined as 1 second—while the agent has to allocate more bandwidth to URLLC slice to guarantee conformity to SLA until the next update, which wastes the bandwidth to some extent and thus leads to the decrease of SE. There are two lessons that we learn from these simulation results: (a) it is non-trivial to optimize multiple conflicting objectives, even when using cutting-edge RL algorithms; (b) shortening the interval between successive bandwidth allocations may improve performance but it also increases computational costs and raises issues of stability due to more drastic changes in demand.

附註(4). 然後當頻率改變的時候，狀態變更的時間比狀態變化還要快，畢竟真是根據時間的。在之後天數 Interval 很久了

V. CONCLUSION

In this paper, we have investigated the combination of deep distributional RL and GAN and proposed GAN-DDQN to learn the optimum solution for demand-aware resource management in network slicing. In particular, we have applied GAN to approximate the action-value distribution, so as to avoid the negative impact of randomness and noise on the reward and grasp much more details therein than the conventional DQN. We have also designed a new update procedure that combines the advantages offered by distributional RL with the training algorithm of WGAN-GP. Furthermore, we have adopted the reward-clipping scheme to enhance the training stability of GAN-DDQN. Besides, we have introduced the dueling structure to the generator (i.e., Dueling GAN-DDQN), so as to separate the state-value distribution and the action advantage function from the action-value distribution and thus avoid the inherent training problem of GAN-DDQN.

→ ex. Dueling DQN 和 Dueling GAN DDQN

Extensive simulations have demonstrated the effectiveness of GAN-DDQN and Dueling GAN-DDQN with superior performance over the classical DQN algorithm. In the future, we will try to further improve the GAN-DDQN mechanism under various scenarios with multiple-metric constraints as well as non-stationary traffic demands.) ↗ ex. Energy Efficiency

APPENDIX THE PROOF OF THEOREM 1

Before the proof of Theorem 1, we give the following lemmas:

Lemma 1: Because $v(\theta, \psi) = 0$ if and only if $(\theta, \psi) = (\xi, 0)$, the unique Nash-equilibrium point of the training objective in Eq. (30) is given by $\theta = \xi, \psi = 0$.

Lemma 2: The distance between the optimal boundaries of Dirac-WGAN-GP on \mathcal{D}_1 and \mathcal{D}_2 is δ .

Proof: Dirac-WGAN-GP consists of a generator with parameter θ and a linear discriminator $D_\psi(x) = \psi \cdot x$, where the generator outputs a Dirac distribution centralized θ (i.e., δ_θ). Whilst the real data distribution is given by a Dirac distribution concentrated at ξ (i.e., δ_ξ). Therefore, the training objective of Dirac-WGAN-GP is given by

$$L(\theta, \psi) = \psi\theta - \xi\psi \quad (30)$$

and the gradient penalty proposed in [42] is given by

$$\begin{aligned} p(\psi) &= \frac{\lambda}{2} \mathbb{E}_{\hat{x}} (\|\nabla_{\hat{x}} D_\psi(\hat{x})\| - 1)^2 \\ &= \frac{\lambda}{2} (|\psi| - 1)^2 \end{aligned} \quad (31)$$

Inspired by [46], we use gradient vector field to analyze convergence, which is defined as follow

$$v(\theta, \psi) := \begin{pmatrix} -\nabla_\theta L(\theta, \psi) \\ \nabla_\psi L(\theta, \psi) \end{pmatrix} \quad (32)$$

For Dirac-WGAN-GP, the corresponding gradient vector field is given by

$$v(\theta, \psi) = \begin{pmatrix} -\psi \\ \theta - \xi + \text{sign}(\psi)\lambda(|\psi| - 1) \end{pmatrix} \quad (33)$$

where $\text{sign}(\cdot)$ denotes the signum function and we have Lemma 1.

Assume that the iteration (θ_k, ψ_k) converges towards the equilibrium point $(\xi, 0)$ but $(\theta_k, \psi_k) \neq (\xi, 0)$ for all $k \in \mathbb{N}$, which implies that $v(\theta_k, \psi_k) \approx 0$ and thus we have

$$-\psi_k \approx \theta_k - \xi + \text{sign}(\psi_k)\lambda(|\psi_k| - 1) \quad (34)$$

in other words,

$$\theta_k \approx -\psi_k + \xi - \text{sign}(\psi_k)\lambda(|\psi_k| - 1) \quad (35)$$

Then, we can get the update amount of parameter θ after the $(k+1)$ th training as follow

$$\begin{aligned} |\theta_{k+1} - \theta_k| &\approx h |-\psi_k + \xi - \text{sign}(\psi_k)\lambda(|\psi_k| - 1) - \theta_k| \\ &\approx h |-(\lambda + 1)\psi_k + (\xi - \theta_k) + \text{sign}(\psi_k)\lambda| \end{aligned} \quad (36)$$

Therefore, we have $\lim_{k \rightarrow \infty} |\theta_{k+1} - \theta_k| = h\lambda$, which shows that Dirac-WGAN-GP cannot converge to the equilibrium

point, and the value of generator's parameter will finally oscillate between $\xi - \frac{h\lambda}{2}$ and $\xi + \frac{h\lambda}{2}$.

Assume that \mathcal{D}_1 and \mathcal{D}_2 are two different real data, which are the Dirac distributions concentrated at ξ_1 and ξ_2 , respectively. Let $\delta = |\xi_1 - \xi_2|$, which indicates the statistic distance between \mathcal{D}_1 and \mathcal{D}_2 . Note that usually $h\lambda$ is two or three orders of magnitude smaller than δ , which implies that the optimal boundaries is rarely overlapped, thus further training is required when the real data varies. With the constant learning rate, it is easy to deduce from Lemma 2 that the larger the δ is, the more training steps are required for the Dirac-WGAN-GP to reach the new optimal boundary. Finally, based on Lemma 1 and 2, we obtain the proof of Theorem 1. ■

REFERENCES

- [1] Y. Hua, R. Li, Z. Zhao, H. Zhang, and X. Chen, "GAN-based deep distributional reinforcement learning for resource management in network slicing," in *Proc. Globecom*, Waikoloa, HI, USA, Dec. 2019.
- [2] K. Katsalis, N. Nikaein, E. Schiller, A. Ksentini, and T. Braun, "Network slices toward 5G communications: Slicing the LTE network," *IEEE Commun. Mag.*, vol. 55, no. 8, pp. 146–154, Aug. 2017.
- [3] R. Li *et al.*, "Intelligent 5G: When cellular networks meet artificial intelligence," *IEEE Wireless Commun.*, vol. 24, no. 5, pp. 175–183, Oct. 2017.
- [4] X. Foukas, G. Patounas, A. Elmokashfi, and M. K. Marina, "Network slicing in 5G: Survey and challenges," *IEEE Commun. Mag.*, vol. 55, no. 5, pp. 94–100, May 2017.
- [5] *Minimum Requirement Related to Technical Performance for IMT-2020 Radio Interface (s)*, document ITU-R M.2410-0, Nov. 2017.
- [6] X. Zhou, R. Li, T. Chen, and H. Zhang, "Network slicing as a service: Enabling enterprises' own software-defined cellular networks," *IEEE Commun. Mag.*, vol. 54, no. 7, pp. 146–153, Jul. 2016.
- [7] X. Li *et al.*, "Network slicing for 5G: Challenges and opportunities," *IEEE Internet Comput.*, vol. 21, no. 5, pp. 20–27, Sep. 2017.
- [8] I. Afolabi, T. Taleb, K. Samdanis, A. Ksentini, and H. Flinck, "Network slicing and softwarization: A survey on principles, enabling technologies, and solutions," *IEEE Commun. Surveys Tuts.*, vol. 20, no. 3, pp. 2429–2453, 3rd Quart., 2018.
- [9] H. Zhang, N. Liu, X. Chu, K. Long, A. Aghvami, and V. C. M. Leung, "Network slicing based 5G and future mobile networks: Mobility, resource management, and challenges," *IEEE Commun. Mag.*, vol. 55, no. 8, pp. 138–145, Aug. 2017.
- [10] J. Ordóñez-Lucena, P. Ameigeiras, D. Lopez, J. J. Ramos-Munoz, J. Lorca, and J. Folgueira, "Network slicing for 5G with SDN/NFV: Concepts, architectures, and challenges," *IEEE Commun. Mag.*, vol. 55, no. 5, pp. 80–87, May 2017.
- [11] I. da Silva *et al.*, "Impact of network slicing on 5G radio access networks," in *Proc. EuCNC*, Athens, Greece, Jun. 2016, pp. 153–157.
- [12] *Study on New Services and Markets Technology Enables, Release 14*, document 3GPP TR 22.981, Mar. 2016.
- [13] R. Li *et al.*, "Deep reinforcement learning for resource management in network slicing," *IEEE Access*, vol. 6, pp. 74429–74441, 2018.
- [14] S. Vassilaras *et al.*, "The algorithmic aspects of network slicing," *IEEE Commun. Mag.*, vol. 55, no. 8, pp. 112–119, Aug. 2017.
- [15] B. Han, J. Lianghai, and H. D. Schotten, "Slice as an evolutionary service: Genetic optimization for inter-slice resource management in 5G networks," *IEEE Access*, vol. 6, pp. 33137–33147, 2018.
- [16] P. L. Vo, M. N. H. Nguyen, T. A. Le, and N. H. Tran, "Slicing the edge: Resource allocation for RAN network slicing," *IEEE Wireless Commun. Lett.*, vol. 7, no. 6, pp. 970–973, Dec. 2018.
- [17] Y. Sun, G. Feng, L. Zhang, M. Yan, S. Qin, and M. A. Imran, "User access control and bandwidth allocation for slice-based 5G-and-beyond radio access networks," in *Proc. ICC*, Shanghai, China, May 2019, pp. 1–6.
- [18] M. Jiang, M. Condoluci, and T. Mahmoodi, "Network slicing management & prioritization in 5G mobile systems," in *Proc. Eur. Wireless Conf.*, Oulu, Finland, May 2016, pp. 1–6.
- [19] S. D'Oro, F. Restuccia, T. Melodia, and S. Palazzo, "Low-complexity distributed radio access network slicing: Algorithms and experimental results," *IEEE/ACM Trans. Netw.*, vol. 26, no. 6, pp. 2815–2828, Dec. 2018.
- [20] Z. Zhou, L. Tan, B. Gu, Y. Zhang, and J. Wu, "Bandwidth slicing in software-defined 5G: A Stackelberg game approach," *IEEE Veh. Technol. Mag.*, vol. 13, no. 2, pp. 102–109, Jun. 2018.
- [21] G. A. Rummery and M. Niranjan, *On-Line Q-Learning Using Connectionist Systems*, vol. 37. Cambridge, U.K.: Univ. of Cambridge, 1994.
- [22] C. J. C. H. Watkins and P. Dayan, "Q-learning," *Mach. Learn.*, vol. 9, nos. 3–4, pp. 279–292, May 1992.
- [23] V. Mnih *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
- [24] M. G. Bellemare, W. Dabney, and R. Munos, "A distributional perspective on reinforcement learning," in *Proc. ICML*, Sydney, NSW, Australia, Aug. 2017, pp. 449–458.
- [25] R. Koenker and K. Hallock, "Quantile regression," *J. Econ. Perspect.*, vol. 15, no. 4, pp. 143–156, 2001.
- [26] W. Dabney, M. Rowland, M. G. Bellemare, and R. Munos, "Distributional reinforcement learning with quantile regression," in *Proc. AAAI*, New Orleans, LA, USA, Feb. 2018, pp. 2892–2901.
- [27] W. Dabney, G. Ostrovski, D. Silver, and R. Munos, "Implicit quantile networks for distributional reinforcement learning," Jun. 2018, *arXiv:1806.06923*. [Online]. Available: <https://arxiv.org/pdf/1806.06923.pdf>
- [28] I. Goodfellow *et al.*, "Generative adversarial nets," in *Proc. NeurIPS*, Montreal, QC, Canada, Dec. 2014, pp. 2672–2680.
- [29] X. Li, J. Fang, W. Cheng, H. Duan, Z. Chen, and H. Li, "Intelligent power control for spectrum sharing in cognitive radios: A deep reinforcement learning approach," *IEEE Access*, vol. 6, pp. 25463–25473, 2018.
- [30] Z. Xu, Y. Wang, J. Tang, J. Wang, and M. C. Gursoy, "A deep reinforcement learning based framework for power-efficient resource allocation in cloud RANs," in *Proc. ICC*, Paris, France, May 2017, pp. 1–6.
- [31] N. Liu *et al.*, "A hierarchical framework of cloud resource allocation and power management using deep reinforcement learning," in *Proc. ICDCS*, Atlanta, GA, USA, vol. Jun. 2017, pp. 372–382.
- [32] Y. He, F. R. Yu, N. Zhao, V. C. M. Leung, and H. Yin, "Software-defined networks with mobile edge computing and caching for smart cities: A big data deep reinforcement learning approach," *IEEE Commun. Mag.*, vol. 55, no. 12, pp. 31–37, Dec. 2017.
- [33] T. Doan, B. Mazoure, and C. Lyle, "GAN Q-learning," May 2018, *arXiv:1805.04874*. [Online]. Available: <https://arxiv.org/pdf/1805.04874.pdf>
- [34] D. Freirich, T. Shimkin, R. Meir, and A. Tamar, "Distributional multi-variate policy evaluation and exploration with the Bellman GAN," *Proc. ICML*, Long Beach, CA, USA, Jun. 2019.
- [35] Z. Wang, N. de Freitas, and M. Lanctot, "Dueling network architectures for deep reinforcement learning," Nov. 2015, *arXiv:1511.06581*. [Online]. Available: <https://arxiv.org/pdf/1511.06581.pdf>
- [36] C. Qi, Y. Hua, R. Li, Z. Zhao, and H. Zhang, "Deep reinforcement learning with discrete normalized advantage functions for resource management in network slicing," *IEEE Commun. Lett.*, vol. 23, no. 8, pp. 1337–1341, Aug. 2019.
- [37] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 1998.
- [38] V. Mnih *et al.*, "Playing Atari with deep reinforcement learning," Dec. 2013, *arXiv:1312.5602*. [Online]. Available: <https://arxiv.org/pdf/1312.5602.pdf>
- [39] J. N. Tsitsiklis and B. Van Roy, "Feature-based methods for large scale dynamic programming," *Mach. Learn.*, vol. 22, no. 1, pp. 59–94, 1996.
- [40] G. Barth-Maron *et al.*, "Distributed distributional deterministic policy gradients," Apr. 2018, *arXiv:1804.08617*. [Online]. Available: <https://arxiv.org/pdf/1804.08617.pdf>
- [41] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein GAN," Jan. 2017, *arXiv:1701.07875*. [Online]. Available: <https://arxiv.org/pdf/1701.07875.pdf>
- [42] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of Wasserstein GANs," in *Proc. NeurIPS*, Long Beach, CA, USA, Dec. 2017, pp. 5767–5777.
- [43] S. Costanzo, I. Fajjari, N. Aitsaadi, and R. Langar, "A network slicing prototype for a flexible cloud radio access network," in *Proc. CCNC*, Las Vegas, NV, USA, Jan. 2018, pp. 1–4.
- [44] S. Costanzo, I. Fajjari, N. Aitsaadi, and R. Langar, "DEMO: SDN-based network slicing in C-RAN," in *Proc. CCNC*, Las Vegas, NV, USA, Jan. 2018, pp. 1–2.
- [45] *Network Slicing Architecture*. Accessed: Oct. 13, 2019. [Online]. Available: <https://tools.ietf.org/id/draft-geng-netslices-architecture-02.html#rfc.references.2>

- [46] L. Mescheder, A. Geiger, and S. Nowozin, "Which training methods for GANs do actually converge?" Jan. 2018, *arXiv:1801.04406*. [Online]. Available: <https://arxiv.org/pdf/1801.04406.pdf>
- [47] *Evolved Universal Terrestrial Radio Access (E-UTRA); Further Advancements for E-UTRA Physical Layer Aspects, Release 9*, document 3GPP TR 36.814, Mar. 2010.
- [48] *Service Requirements for Next Generation New Services Markets, Release 15*, document 3GPP TS 22.261, Mar. 2017.



Yuxiu Hua received the B.S. degree in electronic information engineering from Hangzhou Dianzi University, Hangzhou, China, in June 2016. He is currently pursuing the Ph.D. degree with the College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou. His research interest includes deep learning, reinforcement learning, and network slicing.



Rongpeng Li (S'12–M'17) received the B.E. degree (Hons.) from Xidian University, Xi'an, China, in June 2010, and the Ph.D. degree (Hons.) from Zhejiang University, Hangzhou, China, in June 2015. He was a Research Engineer with the Wireless Communication Laboratory, Huawei Technologies Company, Ltd., Shanghai, China, from August 2015 to September 2016. In November 2016, he returned to academia as a Post-Doctoral Researcher with the College of Computer Science and Technologies, Zhejiang University, which is

sponsored by the National Post-Doctoral Program for Innovative Talents. He is currently an Assistant Professor with the College of Information Science and Electronic Engineering, Zhejiang University. His research interests currently focus on reinforcement learning, data mining, and all broad-sense network problems (e.g., resource management and security). He has authored or coauthored several articles in the related fields. He serves as an Editor of *China Communications*.



Zhifeng Zhao received the bachelor's degree in computer science, the master's degree in communication and information system, and the Ph.D. degree in communication and information system from the PLA University of Science and Technology, Nanjing, China, in 1996, 1999, and 2002, respectively. From 2002 to 2004, he acted as a Post-Doctoral Researcher with the Zhejiang University, Hangzhou, China, where his researches were focused on multimedia next-generation networks (NGN) and soft-switch technology for energy efficiency. From 2005 to 2006, he also acted as a Senior Researcher with the PLA University of Science and Technology, Nanjing, where he performed research and development on advanced energy-efficient wireless router, ad hoc network simulator, and cognitive mesh networking test-bed. From 2006 to 2019, he was an Associate Professor with the College of Information Science and Electronic Engineering, Zhejiang University. He is currently with Zhejiang Lab, Hangzhou, China. His research area includes cognitive radio, wireless multi-hop networks (ad hoc, Mesh, and WSN), wireless multimedia networks, and green communications. He is the Symposium Co-Chair of the ChinaCom 2009 and 2010. He is also the Technical Program Committee (TPC) Co-Chair of the IEEE ISCIT 2010 (10th IEEE International Symposium on Communication and Information Technology).



Xianfu Chen received the Ph.D. degree (Hons.) in signal and information processing from the Department of Information Science and Electronic Engineering (ISEE), Zhejiang University, Hangzhou, China, in March 2012. Since April 2012, he has been with the VTT Technical Research Centre of Finland, Oulu, Finland, where he is currently a Senior Scientist. His research interests cover various aspects of wireless communications and networking, with emphasis on human-level and artificial intelligence for resource awareness in next-generation communication networks. He is serving and served as the Track Co-Chair and a TPC member for a number of IEEE ComSoc flagship conferences. He is also a Vice Chair of the IEEE Special Interest Group on Big Data with Computational Intelligence, a member of which come from more than 15 countries worldwide.



Honggang Zhang was an Honorary Visiting Professor with the University of York, U.K., and an International Chair Professor of Excellence for the Université Européenne de Bretagne and Supèlec, France. He is currently a Full Professor with the College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou, China. He has coauthored and edited two books: *Cognitive Communications: Distributed Artificial Intelligence (DAI), Regulatory Policy and Economics, Implementation* (John Wiley & Sons) and *Green Communications: Theoretical Fundamentals, Algorithms and Applications* (CRC Press), respectively. He is also active in the research on cognitive green communications. He was the leading Guest Editor of the *IEEE Communications Magazine* special issues on Green Communications. He served as the Chair of the Technical Committee on Cognitive Networks of the IEEE Communications Society from 2011 to 2012 and the Series Editor for the *IEEE Communications Magazine* for its Green Communications and Computing Networks Series from 2015 to 2018. He is an Associate Editor-in-Chief of *China Communications*.