

A New Deep-Q-Learning-Based Transmission Scheduling Mechanism for the Cognitive Internet of Things

Jiang Zhu, Yonghui Song, Dingde Jiang, *Member, IEEE*, and Houbing Song[✉], *Senior Member, IEEE*

Abstract—Cognitive networks (CNs) are one of the key enablers for the Internet of Things (IoT), where CNs will play an important role in the future Internet in several application scenarios, such as healthcare, agriculture, environment monitoring, and smart metering. However, the current low packet transmission efficiency of IoT faces a problem of the crowded spectrum for the rapidly increasing popularities of various wireless applications. Hence, the IoT that uses the advantages of cognitive technology, namely the cognitive radio-based IoT (CIoT), is a promising solution for IoT applications. A major challenge in CIoT is the packet transmission efficiency using CNs. Therefore, a new Q-learning-based transmission scheduling mechanism using deep learning for the CIoT is proposed to solve the problem of how to achieve the appropriate strategy to transmit packets of different buffers through multiple channels to maximize the system throughput. A Markov decision process-based model is formulated to describe the state transformation of the system. A relay is used to transmit packets to the sink for the other nodes. To maximize the system utility in different system states, the reinforcement learning method, i.e., the Q learning algorithm, is introduced to help the relay to find the optimal strategy. In addition, the stacked auto-encoders deep learning model is used to establish the mapping between the state and the action to accelerate the solution of the problem. Finally, the experimental results demonstrate that the new action selection method can converge after a certain number of iterations. Compared with other algorithms, the proposed method can better transmit packets with less power consumption and packet loss.

Index Terms—Cognitive networks (CNs), deep learning, Internet of Things (IoT), Markov decision process, Q learning (QL).

Manuscript received March 12, 2017; revised June 30, 2017; accepted September 2, 2017. Date of publication October 4, 2017; date of current version August 9, 2018. This work was supported in part by the National Nature Science Foundation of China under Grant 61102062, Grant 61271260, Grant 61571104, and Grant 61572231, in part by the Key Project of Chinese Ministry of Education under Grant 212145, in part by the Nature Science Foundation of Chongqing Science and Technology Commission under Grant cstc2015jcyjA40050, in part by the Science and Technology Research Project of Chongqing Education Commission under Grant KJ120530, and in part by the Fundamental Research Funds for the Central Universities under Grant N150402003. (Corresponding authors: Dingde Jiang; Houbing Song.)

J. Zhu and Y. Song are with the Chongqing Key Laboratory of Mobile Communications Technology, Chongqing University of Posts and Telecommunications, Chongqing 400065, China (e-mail: zhujiang@cqupt.edu.cn).

D. Jiang is with the School of Astronautics and Aeronautics, University of Electronic Science and Technology of China, Chengdu 611731, China (e-mail: jiangdd@uestc.edu.cn).

H. Song is with the Department of Electrical, Computer, Software, and Systems Engineering, Embry–Riddle Aeronautical University, Daytona Beach, FL 32114 USA (e-mail: h.song@ieee.org).

Digital Object Identifier 10.1109/JIOT.2017.2759728

I. INTRODUCTION

IN THE future, wireless sensor networks are expected to be integrated into the Internet of Things (IoT) [1], [2], where reconfigurable, flexible, and intelligent sensors dynamically join the Internet and use it to collaborate and accomplish their tasks for a wide range of applications in various domains [3]–[8], such as big data applications, IoT, E-commerce, medical device [9], [10], virtual reality and augmented reality, and environment monitoring. The network environment also tends to become increasingly complicated, and the communication resources become increasingly scarce. It is a great challenge to the wireless sensor networks and IoT. Coincidentally, the cognitive network technology can compensate for these deficiencies [11]–[15]. Cognitive nodes are intelligent wireless devices that can sense the environment, observe the network changes, use the knowledge learnt from the previous interaction with the network, and make intelligent decisions to seize the opportunities to transmit. The process of continuously sensing the environment information, exchanging control information, learning information, deciding and executing a strategy in the network can provide the ability of intelligence and adaptability to the wireless sensor networks and the future IoT. Therefore, the cognitive radio technology is a key communication approach for resource-constrained wireless sensor networks and future wireless network [16]–[20]. When cognitive users, i.e., sensors in wireless sensor networks, access the spectrum, to effectively use the network resources and satisfy the throughput demand for multimedia applications, effective mechanisms are required to coordinate the actions of the cognitive users (transmission power control, spectrum access, transmission scheduling, *et al.*) [21]–[24]. With the rapid increase in number of wireless devices in the IoT, more data will be stored in the network nodes. Thus, the method to rapidly forward data with the limited storage space and bandwidth is a great challenge for the current wireless network of IoT.

Currently, many existing literatures (see in [25]–[33]) have studied the problem of network data transmission with unknown environment information in a cross-layer design manner. Among these literatures, [25], [26] have examined the adaptive modulation (AM) algorithm in the data transmission stage, whereas [27] focused on the reliable route discovery to reduce the data online time. Chen *et al.* [28] proposed the quality of service (QoS) awareness scheduler and power

adaptation scheme at both uplink and downlink medium access control (MAC) layers to coordinate the action of the lower layers for resource efficiency. Praveena *et al.* [29] focused on the throughput and fairness. Depending on the preference of the two features, algorithms and methods are proposed to assign or schedule users to prioritize to maximize throughput, maximize fairness or finding the appropriate balance between the two. In [30], a cross layer design of MAC and routing protocols and the topology design are studied to gain more profits for networking. In addition, [31]–[33] modeled the optimization problem of wireless networks as Markov decision process (MDP) to describe the state transformation of the system. However, it is difficult to solve the MDP problem because the MDP has many variables. Therefore, the reinforcement learning method can be introduced to solve it. Lin *et al.* [33], Zheng *et al.* [34], [35], [37], and Duan *et al.* [36] researched the scheduling for different applications or services in distributed wireless networks. The scheduling mechanisms are in cross-layer way, and QoS or quality of experience are involved in the designs of scheduling mechanisms.

For the problem of the curse of dimensionality, a new Q-learning (QL)-based transmission scheduling mechanism using deep learning is proposed for the MDP problem to intelligently make the appropriate strategy. The model formulated in this paper considers the power required for packet transmission in the wireless networks and the packet loss to maximize the system throughput. We modify the action selection by the comprehensive action evaluation method based on the Q value and index value to consider the balance of exploration and exploitation. In addition, when the system scale is large, the state space and action space of the system are also notably large, and it is notably difficult to calculate the optimal action corresponding to the state one by one. Therefore, we use the method of deep learning to construct the mapping between state and action to quickly obtain the strategy.

The contributions of this paper are as follows.

- 1) The MDP model is formulated to describe the problem of transmission scheduling for the cognitive IoT.
- 2) The QL algorithm is modified to learn the system state transition without the system prior information.
- 3) The stacked auto-encoder (SAE) deep learning algorithm is adopted to map the relation between the states and the actions to avoid the massive calculation and storage in the QL phase.

This paper is organized as follows. Sections II and III introduce the system model and MDP problem. Then, Section IV introduces the proposed scheme of this paper, i.e., the deep QL algorithm, and the algorithms are compared in Section V. Section VI provides the scheme simulation and analysis. Finally, we summarize this paper and point out the future direction. Some notations are listed in Table III for the reader's convenience.

II. SYSTEM MODEL

As shown in Fig. 1, a cognitive radio-based IoT (CIoT) that coexists with a licensed system is considered. In the system, a point-to-point transmission from the relay to the sink and M

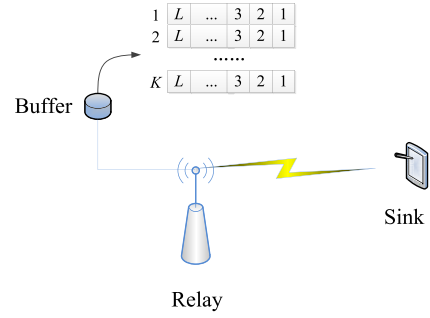


Fig. 1. System model.

frequency domain channels are considered. One relay gathers packets from its K neighbor nodes, and these packets are stored in K buffers with identical length L . The packets that come from K neighbor nodes are assumed to have Poisson distribution with identical arrival rate λ . M channels are independent and identically distributed (i.i.d.). The transmission schedule is decided by the relay. In a frame, the relay selects a channel and transmits packets for a node to the sink. When the channel state is poor, the relay does not transmit packet. When a certain buffer is full, if the relay does not transmit packet for it, then, the packet are lost if packets continue to arrive in the next frame. Therefore, in the packet transmission process, the relay must comprehensively consider the channel state and the buffer states that correspond to the communication pairs and transmission mode.

A. Channel State

In the system model, the time unit is defined as frame T_f , the state of the channel does not change in each frame, and the state transition of the channel occurs between two adjacent states. The channel state can be modeled as a finite-state Markov chain [38]. The signal-to-noise ratio (SNR) is assumed to obey Rayleigh distribution, and the probability density function is expressed as $p(\rho) = 1/\bar{\rho} \exp(-\rho/\bar{\rho})$, where $\rho > 0$; $\bar{\rho} = E(\rho)$ is the average SNR. The threshold of the SNR is $\rho_{\text{snr}} = \{\rho_1, \rho_2, \dots, \rho_{C-1}\}$, and C is the number of channel states. Then, we obtain the channel state space $C \triangleq \{c_0, c_1, \dots, c_{C-1}\}$. Therefore, the probability distribution of channel state is

$$p_C(c_n) = \int_{\rho_n}^{\rho_{n+1}} p(\rho) d\rho. \quad (1)$$

The state transition probability of the channel is

$$p_C(c_n, c_{n+1}) = N(\rho_{n+1})T_f / p_C(c_n), \quad n \in \{1, 2, \dots, N-2\} \quad (2)$$

$$p_C(c_n, c_{n-1}) = N(\rho_n)T_f / p_C(c_n), \quad n \in \{1, 2, \dots, N-1\} \quad (3)$$

where $N(\rho_n) = \sqrt{2\pi\rho_n/\bar{\rho}} f_d \exp(-\rho/\bar{\rho})$; f_d is the maximum Doppler shift. Since M channels are i.i.d., the state transition probability of M channels is $p_C(c, c') = \prod_{m=1}^M p_{C,m}(c_i, c_{i+1})$.

B. Buffer State

In each system frame T_f , the arriving packet obeys Poisson distribution with arrival rate λ , which is expressed as

$p_{d_i}(d_i) = \exp(-\lambda_{d_i} T_f) (\lambda_{d_i} T_f)^{d_i} / d_i!$, where d_i is the number of arriving packets in each frame. At the beginning of frame i , to buffer k , the existing buffer length is $l_{i,k}$. If the number of arriving packets is $d_{i,k}$ and the number of transmitted packets is $t_{i,k}$, then the new buffer length is

$$l_{k,i+1} = \min(l_{i,k} + d_{i,k} - t_{i,k}, L). \quad (4)$$

If the state transition probability of buffer k is $p_{l_k}(l_k, l'_k)$, the state transition probability of the K buffers is $p_l(l, l') = \prod_{k=1}^K p_{l_k}(l_k, l'_{i,k})$, where $a_{i,k}$ is the number of packets transmitted for buffer k during frame i .

C. Transmission Power

To truly improve the transmission efficiency, the AM [26] method is used to adjust the transmission power and rate. We use $j \in \{0, 1, 2, \dots, J\}$ to indicate the selected mode. 0 and 1 are corresponding to no transmission and BPSK transmission, respectively, and $j \geq 2$ is corresponding to 2^j -QAM transmission. Given transmission rate, power, and channel state, the bit error rate (BER) can be estimated. Assuming ideal coherent phase detection, BER bounds are given by [39]

If $j = 1$, then

$$p_{\text{BER}}(c_i, j) \leq 0.5 \operatorname{erfc} \left(\sqrt{\rho_i P(c_i, j) / W N_0} \right) \quad (5)$$

If $j > 1$, i.e., $j = 2, 3, \dots$, then

$$p_{\text{BER}}(c_i, j) \leq 0.2 \exp(-1.6 \rho_i P(c_i, j) / W N_0 (2^j - 1)) \quad (6)$$

where $W N_0$ is the noise power. The BER inequalities above give a pessimistic minimum power $P(c_i, j)$ to achieve a specified BER for channel state c_i and selected mode j .

III. MDP ANALYSIS OF TRANSMISSION

The system contains two state objects: 1) buffer state and 2) channel state. The system operation is a process of state transition. The next system state is obtained by selecting and implementing a certain action at the current system state. Therefore, the state of the system in the next frame is only related to the current state and action. Thus, we model the transmission scheduling problem as Markov decision process.

A. Action Set

When state transition occurs, the relay must choose an action according to current state. The possible action of relay can be defined as $a_i \in A = \{a_{m,k,j}\}$, where $m \in \{1, 2, \dots, M\}$, $k \in \{1, 2, \dots, K\}$, and $j \in \{0, 1, 2, \dots, J\}$. $a_i = a_{m,k,j}$ indicates that at the beginning of frame i , the relay selects channel m to transmit $a_{m,k,j}$ packets by selected mode j for buffer k .

B. State Transition Probability

The system state is the combination of the buffer state and channel state, i.e., $S \triangleq B \otimes C$. If the buffer length is L , the number of buffer states of a single buffer is $B = L + 1$. The number of channels is M .

To one buffer and one channel, the state transition probability is $p_l(l_i, l_{i+1} | a_i) \times p_c(c_i, c_{i+1})$. Therefore, the entire system state transition probability expression is

$$p_s(s_i, s_{i+1} | a_i) = \prod_{k=1}^K p_{l_k}(l_i, l_{i+1} | a_i) \times \prod_{m=1}^M p_{c_m}(c_i, c_{i+1}). \quad (7)$$

C. Utility

The goal of this paper is to maximizing the system utility. If the coding rate is V , the throughput (bits per symbol) under different transmission mode is $V \times j$. Thus, at system $s_i = \{l_i, c_i\} \in S$, the benefit of system after taking action a_i is

$$B(s_i, a_i) = \sum_{k=1}^K V \times j \quad (8)$$

where $R(s_i, a_i)$ is the throughputs when the system state is s_i , and action a_i is selected in frame i . It is mentioned in above section that adaptive transmission scheme based on M-QAM is employed for each channel. When $a_{m,k,j}$ packets are transmitted through channel m , corresponding transmission mode j is employed. Therefore, there is a mapping $\varphi(\cdot)$, that $\varphi(j) = a_{m,k,j}$. For simplification, we assume that $V \times j = a_{m,k,j}$.

The pressure value of buffer k is defined as $f_{k,i} = \exp(\theta \times l_{i,k})$, where θ is the pressure coefficient, which represents the number of packets in the buffer. When the number of packets in buffer increases, the arriving packets in the next frame may be lost because of the small space. Therefore, a smaller pressure value corresponds to less packets loss. Thus, the pressure value is inversely proportional to the performance of the system. In addition, if we have the constraint of the BER, we can obtain the minimum power transmission power $p_{s_i}(s_i, a_i)$ [40], which is also inversely proportional to the performance of the system.

Consequently, the system cost is the combination of buffer pressure and power consumption

$$C(s_i, a_i) = \left(\sum_{i=1}^K f_{k,i} \right) \times p_{s_i}(s_i, a_i). \quad (9)$$

The system utility O_i is proportional to the number of transmitted packets in each frame and inversely proportional to the buffer pressure and power consumption. Hence, the system utility can be described as

$$O_i = O(s_i, a_i) = B(s_i, a_i) / C(s_i, a_i). \quad (10)$$

IV. DEEP-Q-LEARNING-BASED OPTIMAL ACTION ACQUISITION SCHEME

In the transmission process, the relay obtains the environment state information by learning to guide its action. The state and action of the system are discrete, and the state of the system discontinuously changes when an action is executed in a frame. Accordingly, to solve the MDP problem, the method of reinforcement learning is used to guide the action of node [41]–[44]. The QL algorithm learns the environment state information and obtains the optimal action. The QL algorithm is a type of gradual optimization process, so it is difficult

to achieve fast convergence in action selection. The artificial neural network can compensate the limitation for its generalization and function approximation ability. In addition, the deep learning method has proven good performance in many fields [45]–[47]. Thus, we use the deep learning method to establish a mapping between states and actions [48]–[52].

A. Q Learning Algorithm

In the process of the QL algorithm, the agent finds the optimal action by continuous interaction with the environment in a constant trial-and-error manner. The optimal action is related to immediate reward and considers the reward of future n steps. We use $V^\pi(s_i)$ to represent the entire reward it in QL algorithm under a policy π

$$V^\pi(s_i) = r_i + \gamma r_{i+1} + \gamma^2 r_{i+2} + \dots \quad (11)$$

In the QL algorithm, the Q value, which is the evaluation of the state and action, is the combination of immediate and discounted reward and can be formulated as

$$Q(s_i, a_i) \leftarrow r_i + \gamma V^\pi(s_{i+1}) \quad (12)$$

where $\gamma (0 < \gamma < 1)$ is the discount coefficient, which indicates the effect of future reward on the current action. The learning goal of QL is to maximize the total utility. Thus, in (12), we replace r_i and $V^\pi(s_{i+1})$ with O_i and $\max_{a_{i+1} \in A} Q(s_{i+1}, a_{i+1})$, respectively,

$$Q(s_i, a_i) \leftarrow O_i + \gamma \max_{a_{i+1} \in A} Q(s_{i+1}, a_{i+1}) \quad (13)$$

where A is the action sets.

In the learning stage, how to balance the exploration and exploitation of the action set is the key problem in QL. In particular, when the state of the system is large, how to effectively choose the action will directly affect the algorithm convergence and system performance. Therefore, to achieve the optimal action, we add a modified index value to quickly find the best action. The index value can reflect the fluctuations of the rewards and timely adjust the explore range to reduce the unnecessary selection cost

$$a_i \leftarrow \arg \max_a (Q(s_i, a) + \text{Index}(s_i, a)) \quad (14)$$

where Q is the evaluation value to the current state and action. Based on the Q value, $\text{Index}(s_i, a)$ is introduced to obtain the optimal potential action. Its expression is

$$\text{Index}(s_i, a) = C_p \sqrt{2 \ln i \times \min\{1/4, V_a(i)\} / T_a(i)} \quad (15)$$

where C_p is a constant greater than 0 [53]. $T_a(i)$ is the number of times action a has been selected after i frames. $V_a(i)$ is the bias factor, which contains the utility value variance $\sigma_a^2(i)$ of the action to reflect its volatility

$$\sigma_a^2(i) = \sum_{k=1}^{T_a(i)} O^2(s_k, a) / T_a(i) - \bar{O}^2(s_{T_a(i)}, a) \quad (16)$$

$$V_a(i) = \sigma_a^2(i) + \sqrt{2 \ln i / T_a(i)}. \quad (17)$$

On one hand, the action selection method based on the action index considers the system utility of the current action

Algorithm 1 Modified QL Algorithm

1. Initialize action visiting number $T_i(n) = 0$.
2. Initialize state-action value $Q(s_i, a_i) = 0$ and state-action look-up table.
3. for $episode1=1, I_1$ do
4. Initialize action vector $a = \{a_1, a_2, \dots\}$.
5. To the current state s_i .
6. for $episode2=1, I_{te}$ do
7. if $episode2=1$
8. Select a random action a_i .
9. end if
10. if $episode2 > 1$
11. $\text{Index}(s_i, a) \leftarrow C_p \sqrt{\frac{2 \ln i}{T_a(i)} \min\{1/4, V_a(i)\}}$
12. Select action according to the following formula
$$a_i \leftarrow \max_a (Q(s_i, a) + \text{Index}(s_i, a)).$$
13. end if
14. Execute a_i and obtain O_i and turn into state s_{i+1} .
15. Calculate $\alpha \leftarrow 1/(1 + T_a(i))$.
16. Update $Q(s_i, a_i)$.
$$Q_{i+1}(s_i, a_i) \leftarrow \begin{cases} (1 - \alpha)Q_i(s_i, a_i) + \alpha(O_i + \gamma \max_{a_{i+1}} Q_i(s_{i+1}, a_{i+1})) & \text{if } s = s_i \text{ and } a = a_i, \\ Q_i(s_i, a_i) & \text{otherwise.} \end{cases}$$
17. Update state-action look-up table.
18. end for
19. end for

and gradually considers the action with a larger effect, which reflects the characteristic of exploitation of the system. On the other hand, with the ongoing iterative process, if a certain action is not selected or the selected number is notably small, then it is biased to select the action in the next iteration, which reflects the characteristic of exploration.

After determining the execution action, the relay performs action a_i , calculates utility value O , and updates the Q value according to formula

$$Q_{i+1}(s_i, a_i) = \begin{cases} (1 - \alpha)Q_i(s_i, a_i) + \alpha(O_i + \gamma \max_{a_{i+1}} Q_i(s_{i+1}, a_{i+1})), & \text{if } s = s_i \text{ and } a = a_i \\ Q_i(s_i, a_i), & \text{otherwise} \end{cases} \quad (18)$$

where $\alpha (0 < \alpha \leq 1)$ is the learning rate of the state action and is calculated by $\alpha = 1/(1 + T_a(i))$.

The specific implementation process of the modified QL algorithm is shown in Algorithm 1.

B. Convergence Analysis of Algorithm

In the convergence analysis process of algorithm, the optimal Q value is $Q^*(s_i, a_i)$.

Theorem 1: The value of the system utility function O_i , which is defined by formula (10) is bounded.

Proof: Formula (10) consists of two parts: 1) benefit and 2) cost. The expression of benefit is $\sum_{k=1}^K V \times j$, which indicates the number of transmitted packets in a frame, is a finite value. The denominator is the cost function (9) which is the combination of buffer pressure and power consumption. The number of packets in the buffer and power consumption are finite. Consequently, the system utility is bounded. ■

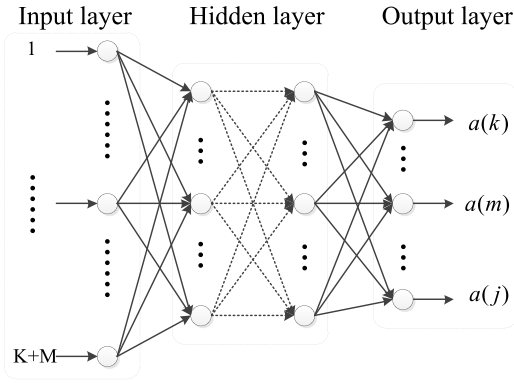


Fig. 2. Structure of the SAE model.

Theorem 2: Learning rate satisfies that $0 < \alpha \leq 1$, and

$$\sum_{i=1}^{\infty} \alpha = \infty, \sum_{i=1}^{\infty} \alpha^2 < \infty \quad \forall s, a. \quad (19)$$

Proof: Since $\alpha = 1/(1 + T_a(i))$ and $T_a(i)$ is the number of times action a has been selected after i frames, $0 < \alpha \leq 1$.

According to the definition of $T_a(i)$, $1/(1 + T_a(i)) < 1/i$. It is clear that $\sum_{i=1}^{\infty} (1/i)^2 < \infty$, therefore, $\sum_{i=1}^{\infty} \alpha^2 < \infty$.

If each state-action pair is visited infinitely, it is obviously that

$$\sum_{i=1}^{\infty} \alpha = \sum_{i=1}^{\infty} 1/(1 + T_a(i)) > \sum_{i=1}^{\infty} 1/i. \quad (20)$$

Since $\sum_{i=1}^{\infty} 1/i = \infty$, $\sum_{i=1}^{\infty} \alpha = \infty$. ■

Theorem 3 [54]: Given bounded utility function O_i , learning rate $0 < \alpha \leq 1$, and $\sum_{i=1}^{\infty} \alpha = \infty$, $\sum_{i=1}^{\infty} \alpha^2 < \infty$, $\forall s, a$, then

$$\lim_{i \rightarrow \infty} Q_i(s, a) = Q^*(s, a) \quad \forall s, a \quad (21)$$

with probability 1. $Q^*(s, a)$ is Q value under the optimal policy π^* .

C. Deep Action Mapping Network

The SAE model of the deep learning is used to build the relations of states and actions. The structure of the model is shown in Fig. 2.

The input layer of the model represents the state information of the system, and the number of neurons in the layer is $K+M$. The input vector can be expressed as $\text{Input} = [l_1, \dots, l_k, \dots, l_K, c_1, \dots, c_m, \dots, c_M, \dots]$. The output layer, which represents the selection action information, can be expressed as $\text{Output} = [a(k), a(m), a(j)]$. It consists of the selected channel m , transmission mode j and buffer k . Its neuron number is $K+M+J$. The hidden layer consists of multiple layers. According to the previous experience, the number of neurons in hidden layer is

$$n_h = \sqrt{n_i + n_o} + \text{Con} \quad (22)$$

where n_i is the number of input layers, n_o is the number of output layers, and n_h is the number of hidden layers. Con is a constant limited in [1, 10].

The SAE model uses the logistic sigmoid function as the transfer function in the encoding and decoding process. The cost function $L(X)$ is defined as follows:

$$L(X) = \arg \min_{x \in (0,1)} \sum_{i=1}^N \|x^i - f(x^i)\|^2 / 2 + \mu J_{W^l} \quad (23)$$

where $J_{W^l} = \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (W_{ji}^l)^2$ is the weight decay term, whose function is to reduce the range of the weight to prevent over-fitting in the training section. The rules of weight and bias vector updating are as follows:

$$W(k+1) = W(k) - \beta \times \partial L / \partial W \quad (24)$$

$$b(k+1) = b(k) - \beta \times \partial L / \partial b \quad (25)$$

where β is the learning rate. The partial derivative of the cost function to the weight and bias parameters are calculated in

$$\begin{aligned} \partial L / \partial W &= -(y'_{m+1,i} - y_{m+1,i}) \times f'_m \times y_{m,i} + \mu W_{j,i}^m \\ &= -\delta_j^{m+1} \cdot y_{m,i} + \mu W_{j,i}^m \end{aligned} \quad (26)$$

$$\begin{aligned} \partial L / \partial b &= -(y'_{m+1,i} - y_{m+1,i}) \times f'_m \\ &= -\delta_j^{m+1} \end{aligned} \quad (27)$$

where m is the m th layer of SAEs; $y_{m,i}$ is the expected values; $y'_{m,i}$ is the current values; and f is the sigmoid function. The back propagation of the residual error δ is

$$\delta_i^m = \delta_j^{m+1} \times W_{j,i}^m \times f'_m. \quad (28)$$

D. Algorithm Description

The Q-learning-based SAEs algorithm flow chart is shown in Fig. 3. In the interaction process with the environment, we will not train the SAE model in the first certain time because we have no optimal actions. With time, increasingly more optimal action will be found and stored in the state-action look-up table. Then, the SAE model will be trained according to the sufficiently optimal state action information. Therefore, when the system transfers into the hidden state, we use the trained SAE network to map between the state and the action and find the optimal action. Then, we execute the new action and update the look-up table. In the follow-up time, if the system transfers to the learned state, the relay will query the state-action table to obtain the executable action.

V. ALGORITHM COMPARISONS

Now, we compare the complexity of each algorithm, as shown in Table I. In the system, the number of buffers is K , the buffer length is L , the number of channels is M , the number of channel states is C , and the number of possible transmission modes is J . Accordingly, the number of system states is $S = (L+1)^K \times C^M$. To each system state, the possible action number is $A = K \times M \times (J+1)$. The system scale is $D = S \times A$.

To verify the performance of the algorithm, we compare it with the strategy iteration (SI) algorithm [33], W learning (WL) algorithm [32], and random selection (RS) algorithm.

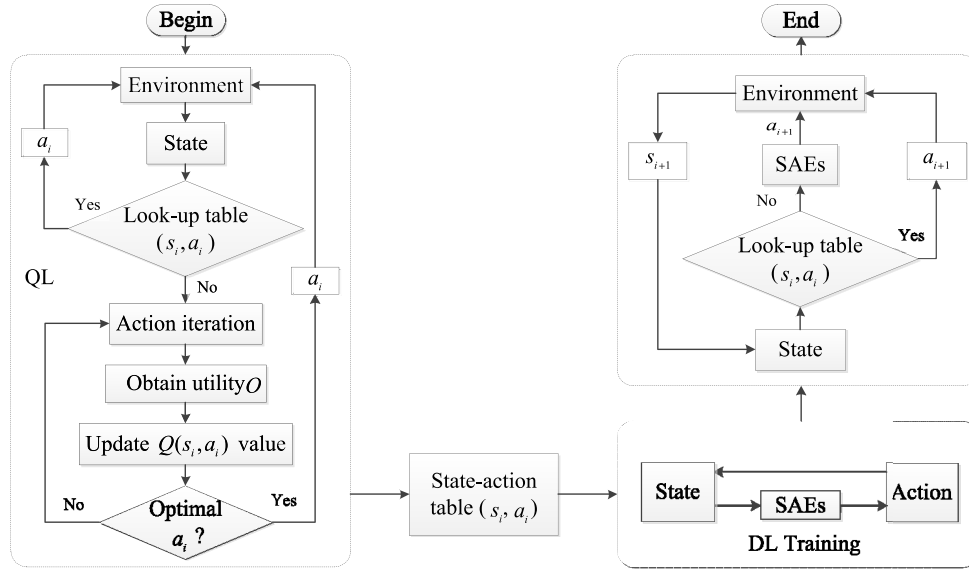


Fig. 3. Flow chart of the proposed scheme.

TABLE I
ALGORITHM COMPLEXITY COMPARISON

Algorithm	Exponential Operation	Multiplication and Division	Addition and Subtraction	Comparison Operation
SI	0	$D + S$	D	$S \times A$
QL	A	$4A$	$2A$	A
WL	$2A$	$3A$	$3A$	A
RS	0	1	0	0

A. Strategy Iteration Algorithm

We can use the SI algorithm to obtain the optimal action as described in (29)

$$V_{n+1} = \max_{a \in A} [r(a) + \gamma p(a) V_n] = r(a_n) + \gamma p(a_n) V_n. \quad (29)$$

The SI algorithm must know the system state transition probability information. However, when the system scale is notably massive, the number of linear equations is equal to the system states [i.e., $S = (L + 1)^K \times C^M$]. Thus, the calculation would be notably massive. It is notably hard to solve because of the problem of dimensional disaster. Therefore, the SI algorithm is not practical in wireless networks.

B. W Learning Algorithm

In the WL method, the QL method is used to obtain the Q value; then, the obtained value is used for WL. The W value represents the difference between the expected return and the actual return

$$W_{i+1}(s_i) = (1 - \alpha)W_i(s_i) + \alpha \left(Q_i(s_i, a_i) - \left(r_i + \gamma \max_{a_i \in A} Q_i(s_{i+1}, a_i) \right) \right). \quad (30)$$

C. Random Selection Algorithm

The RS algorithm is that each system frame randomly selects an action to execute.

In the iteration process, the SI algorithm prefers to calculate all states of the system in one iteration time, whereas the proposed algorithm in this paper prefers to calculate according to the current state. The difference of the RS algorithm is the method of action selection. The RS algorithm randomly selects an action to perform with no extra calculation, so its computational overhead is small. For the QL phrase, the calculation overhead of the algorithm in this paper is equivalent to that of the WL algorithm and lower than that of the SI algorithm.

For the storage overhead, the SI algorithm must calculate the system state transition probability, and the storage overhead is $3(L + 1)^K \times C^M$. The storage in the strategy calculation phase is $(L + 1)^K \times C^M$, whereas the proposed method does not need the prior information storage overhead, and the QL storage overhead is $K \times M \times (J + 1)$, which is much less than SI. The deep learning model requires some storage overhead. However, the model structure in this paper is relatively simple. Therefore, in general, the storage overhead of the proposed method is not larger than SI.

VI. SIMULATION EXPERIMENT AND ANALYSIS

A. Simulation Setting

In the simulation, we define the parameters of the system as follows. The number of wireless communication pairs is $K = 3$, so the number of buffers is $K = 3$. The length of each buffer is $L = 5$. The number of channels is $M = 2$. The number of channel states is $C = 4$. The possible transmission modes are no transmission, QPSK, 4-QAM, 8QAM, and 16QAM, therefore $j \in \{0, 2, 4, 8, 16\}$ and $J = 4$. Coding rate V is 2. For the problem, we compared the three described algorithms with the algorithm of this paper. The simulation parameters were set as shown in Table II. We use the SAE model with two hidden layers. Further performance comparisons are shown in part B.

TABLE II
SIMULATION PARAMETER SETTING

Parameters	Value/Description
Threshold of SNR /dB	$snr = [-6.28, -1.28, 1.28]$
Doppler frequency shift /Hz	$f_d = 50$
Frame length /s	$T_f = 2 \times 10^{-3}$
Slot number I_1	$I_1 = 5 \times 10^3$
Slot number I_2	$I_2 = 1 \times 10^3$
Noise power WN_0/W	1×10^{-3}
Buffer pressure coefficient	$\theta = 0.5$
Arrival rate	$\lambda = [0.1, \dots, 0.9]$
BER constraint	$BER \leq 10^{-3}$
Discount coefficient	$\gamma = 0.9$
Index weight	$C_p = 1/\sqrt{2}$
MQL Learning rate	$\alpha \in (0, 1]$
SAEs hidden neuron number	[8,15,15,9]
Weight	$\mu = 3 \times 10^{-3}$
SAEs Learning rate	$\beta = 1 \times 10^{-2}$
Training error accuracy	$rate = 1 \times 10^{-5}$

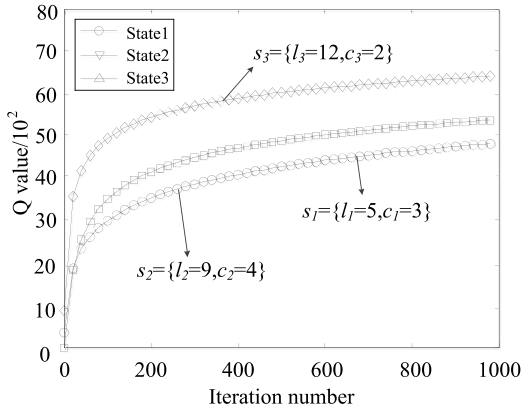


Fig. 4. Q value changing curve.

The SAEs with many hidden layers cannot be fully trained when the data quantity is small. More errors are also introduced because there are more hidden layers. Therefore, the number of hidden layers is selected according to the simulation comparison.

B. Performance Comparison

In slots I_1 , the system learns the optimal actions that correspond to the system states. The optimal state action information is stored in the look-up table. Then, we train the SAE model with the obtained information and use the trained model to map the state and action in the next I_2 phrase. Fig. 4 shows the Q value changing curve in the process of action iteration selection using the QL algorithm.

Fig. 4 shows that when $\lambda = 0.1$, in states $s_1 = \{l_1 = 5, c_1 = 3\}$, $s_2 = \{l_2 = 9, c_2 = 4\}$, and $s_3 = \{l_3 = 12, c_3 = 2\}$, the curve in the modified QL algorithm changes value. From Fig. 4, we observe that the three different curves converge to different values because in different system states, the existing

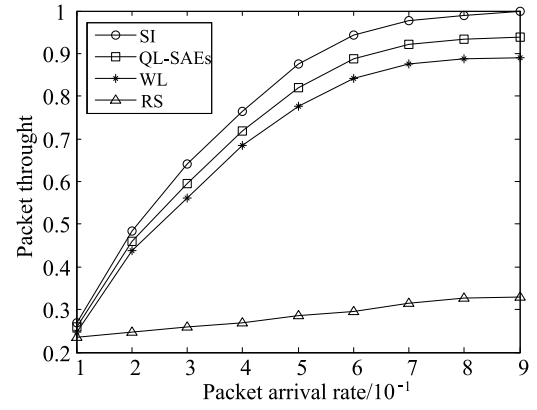


Fig. 5. Normalized throughput comparison of the algorithms.

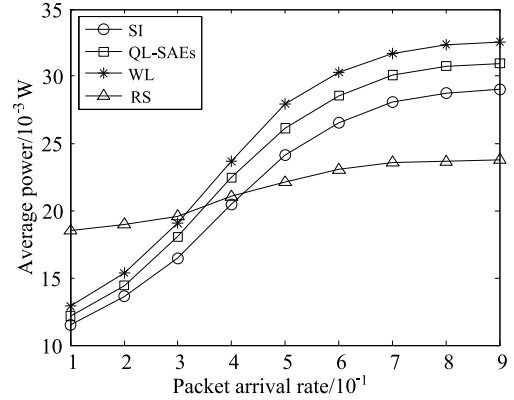


Fig. 6. Average power comparison of the algorithms.

number of packets in buffer, amount of coming packets, and selection of the transmission mode are different. Therefore, the three values converge to different values, and the convergence of the algorithm is illustrated.

In the following, we compare its performance with other algorithms and use QL-SAEs to name our algorithm.

Fig. 5 shows that the system normalized throughput in one frame varies with different packet arrival rates. The QL-SAEs algorithm has a less throughput than the SI algorithm but more than the other two algorithms. This graph shows that when the packet arrival rate is small, SI, QL-SAEs, and WL algorithms have almost equivalent throughputs. Because the buffer pressure is relatively small, in the learning process, the energy consumption is small (Fig. 6) in exchange for the packet throughput. When the packet arrival rate increases, the pressure of the buffer gradually increases. Then, the average amount of packet throughput in each frame gradually increases. Relatively, the arrival rate has less effect on the RS algorithm, and it can obtain a larger throughput only when there are sufficient packets in buffer and better transmission mode is selected.

Fig. 6 shows the average power contrast of each algorithm at different packet arrival rates. When the number of packets is large, the buffer pressure is high, which can force the relay to select better transmission mode and channel to transmit more packets to reduce the buffer pressure. Ultimately, more power

TABLE III
NOTATIONS IN THIS PAPER

Notation	Description	Notation	Description
T_f	Length of frame	$p_c(c, c')$	Channel state transition probability
ρ	Ratio of Signal-to-Noise	f_d	Maximum Doppler shift
C	Channel state space	WN_0	Noise power
$p_c(c_n)$	Channel state probability	λ	Data arrival rate
S	System state space	L	Length of a signal buffer
$p_l(l, l')$	Buffer state transition probability	d	Number of arriving packets
a_i	The selected action	t	Number of transmitted packets
p_{BER}	Bit error rate restriction	$P(c_i, j)$	Power consumption
j	The selected modulation	A	Action set
B	Buffer space	$p_s(s_i, s_{i+1} a_i)$	System transition probability
V	Coding rate	K	Number of buffers
$B(s_i, a_i)$	System benefit	$f_{k,i}$	Pressure of buffer k
θ	Buffer pressure coefficient	$C(s_i, a_i)$	System cost
O_i	System utility	$V^\pi(s_i)$	System long term reward
$Q(s_i, a_i)$	Evaluation of the state and action	γ	Discount coefficient
$Index(s_i, a)$	Action index	$T(n)$	Selection number of the action
$V_a(i)$	Bias factor	$\sigma_i^2(n)$	Action utility value variance
α	Learning rate	$Output$	Output vector of SAEs
$Input$	Input vector of SAEs	n_i	Neuron number in SAEs input layer
M	Number of available channels	n_o	Neuron number in SAEs output layer
Con	A constant limited in [1,10]	n_k	Neuron number in SAEs hidden layer
$L(X)$	Cost function	x^i	Input data
J_{w^l}	Weight decay term	$f(x^i)$	Output data
μ	Weight decay term coefficient	W_{ji}^l	Weight
β	Learning rate	δ	Residual error

is consumed. The energy consumption of the three algorithms rapidly increased at first and subsequently slowly increased because with the increase of packet in the buffer, more packets can be transmitted, and more power is consumed. Therefore, the curve rapidly increases. The buffer will not have much effect on the action when the buffer reaches its limitation. Finally, the power curve tends to gently increase. The power of the RS algorithm is relatively stable because the packet arrival rate basically does not affect the transmission mode selection, whereas the other three algorithms are greatly affected by λ .

The other three algorithms consider the buffer pressure and transmission power, but the curves are above the curve of the RS algorithm.

When the system buffer space is small, in the next frame, more packet will be possibly lost if more packets are arriving. As shown in Fig. 7, with the increase in packet arrival rate, the growth of the packet loss curves is approximately linear. Because the RS algorithm does not consider the power consumption and buffer pressure during the selection process, its packet loss is large compared to other algorithms.

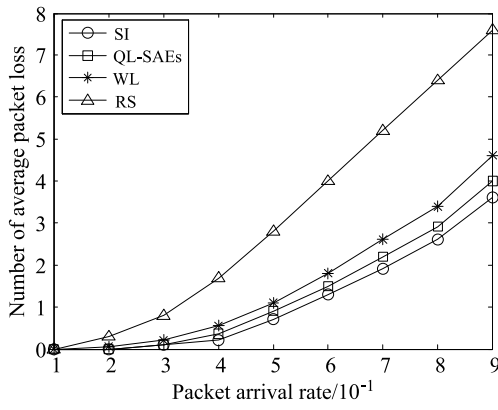


Fig. 7. Number of average packet loss of the algorithms.

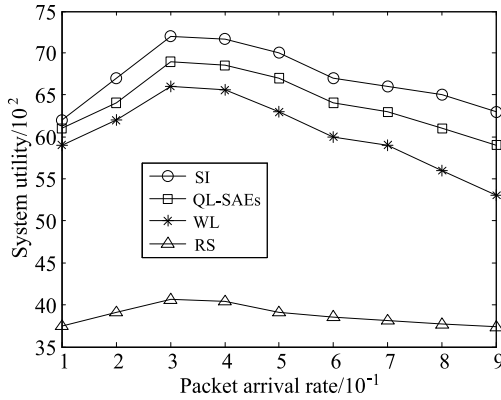


Fig. 8. Average system utility of the algorithms.

The average utility of the system is shown in Fig. 8. The graph shows that the utilities of SI, QL-SAEs, and WL algorithm are higher than RS. Although QL-SAEs have a lower system utility value than SI, it is better than the WL algorithm. To constrain the buffer space, the utility of the system is not notably large when the arrival rate of the packet is too large or too small. When λ is small, the system can select the appropriate action to enhance the utility value. However, when the amount of packet is large, although the system is attempting to transmit packets, it cannot completely transmit all packets. In addition, the transmission power consumption is notably large if there are many packets.

VII. CONCLUSION

In this paper, a relay for transmission packets to the sink of the other nodes is considered in the CIoT. To solve the problem of transmission scheduling in the CIoT, a new QL-based transmission scheduling mechanism using deep learning has been proposed to achieve the appropriate strategy with multiple channels for the cognitive node. In this paper, the reinforcement learning, which is based on the joint action selection criteria of the index value and Q value to balance the action equilibrium problem in exploration and exploitation, is used to solve the MDP problem. Ultimately, it realizes the long-term maximum utility of the system. In addition, the SEA deep learning model is introduced to map between the state and

the action. Although the proposed scheme has a lower performance than SI, the algorithm complexity is strongly reduced. The proposed algorithm can work without priori information, which is suitable for practical scenarios.

The research in this paper is based on the case of one relay. However, if there are more than one relay scenario, the method to make many relays cooperatively or competitively work can be further studied.

REFERENCES

- [1] L. Mainetti, L. Patrono, and A. Vilei, "Evolution of wireless sensor networks towards the Internet of Things: A survey," in *Proc. IEEE Int. Conf. Soft. Telecommun. Comput. Netw.*, Split, Croatia, Sep. 2011, pp. 15–17.
- [2] N. Khalil, M. R. Abid, D. Benhaddou, and M. Gerndt, "Wireless sensors networks for Internet of Things," in *Proc. IEEE 9th Int. Conf. Intell. Sensors Sensor Netw. Inf. Process.*, Singapore, 2014, pp. 1–6.
- [3] S. Jeschke, C. Brecher, H. Song, and D. B. Rawat, Eds., *Industrial Internet of Things: Cyber-Manufacturing Systems*. Cham, Switzerland, Springer, 2016.
- [4] Y. Lin, J. Yang, Z. Lv, W. Wei, and H. Song, "A self-assessment stereo capture model applicable to the Internet of Things," *Sensors*, vol. 15, no. 8, pp. 20925–20944, 2015.
- [5] J. Yang, S. He, Y. Lin, and Z. Lv, "Multimedia cloud transmission and storage system based on Internet of Things," *Multimedia Tools Appl.*, vol. 76, no. 17, pp. 17735–17750, 2015.
- [6] H. A. Maw, H. Xiao, B. Christianson, and J. A. Malcolm, "BTG-AC: Break-the-glass access control model for medical data in wireless sensor networks," *IEEE J. Biomed. Health Inform.*, vol. 20, no. 3, pp. 763–774, May 2016.
- [7] B. Yuan, C. Fu, and D. Chen, "Building a large scale wireless sensor network for the industrial environment," in *Proc. IEEE Int. Conf. Embedded Real Time Comput. Syst. Appl.*, Daegu, South Korea, 2016, pp. 96–96.
- [8] Z. Ma and X. Pan, "Agricultural environment information collection system based on wireless sensor network," in *Proc. IEEE Glob. High Tech Congr. Electron.*, Shenzhen, China, 2012, pp. 24–28.
- [9] M. Chen, J. Yang, Y. Hao, S. Mao, and K. Hwang, "A 5G cognitive system for healthcare," *Big Data Cogn. Comput.*, vol. 1, no. 1, pp. 2–16, 2017, doi: [10.3390/bdcc1010002](https://doi.org/10.3390/bdcc1010002).
- [10] M. Chen, Y. Hao, K. Hwang, L. Wang, and L. Wang, "Disease prediction by machine learning over big data from health-care communities," *IEEE Access*, vol. 5, pp. 8869–8879, 2017, doi: [10.1109/ACCESS.2017.2694446](https://doi.org/10.1109/ACCESS.2017.2694446).
- [11] J. Tervonen, K. Mikhaylov, S. Pieskä, J. Jämsä, and M. Heikkilä, "Cognitive Internet-of-Things solutions enabled by wireless sensor and actuator networks," in *Proc. IEEE CogInfoCom.*, Vietri sul Mare, Italy, 2014, pp. 97–102.
- [12] J. Llore et al., Eds., *Cognitive Networks: Applications and Deployments*. Boca Raton, FL, USA: CRC Press, 2014.
- [13] Q. Wu et al., "Cognitive Internet of Things: A new paradigm beyond connection," *IEEE Internet Things J.*, vol. 1, no. 2, pp. 129–143, Apr. 2014.
- [14] A. A. Khan, M. H. Rehmani, and A. Rachedi, "When cognitive radio meets the Internet of Things?" in *Proc. Int. Wireless Commun. Mobile Comput. Conf.*, Paphos, Cyprus, 2016, pp. 469–474.
- [15] M. A. Shah, S. Zhang, and C. Maple, "Cognitive radio networks for Internet of Things: Applications, challenges and future," in *Proc. Int. Conf. Autom. Comput.*, London, U.K., 2013, pp. 1–6.
- [16] T. Haustein et al., "Cognitive wireless communications—A paradigm shift in dealing with radio resources as a prerequisite for the wireless network of the future—An overview on the topic of cognitive wireless technologies," *Frequenz*, vol. 70, nos. 7–8, pp. 281–288, 2016.
- [17] D. T. Otermat, I. Kostanic, and C. E. Otero, "Analysis of the FM radio spectrum for secondary licensing of low-power short-range cognitive Internet of Things devices," *IEEE Access*, vol. 4, pp. 6681–6691, 2016.
- [18] A. Somov, C. Dupont, and R. Giaffreda, "Supporting smart-city mobility with cognitive Internet of Things," in *Proc. Future Netw. Mobile Summit.*, Lisbon, Portugal, 2013, pp. 1–10.
- [19] M. Nitti, M. Murrioni, M. Fadda, and L. Atzori, "Exploiting social Internet of Things features in cognitive radio," *IEEE Access*, vol. 34, pp. 9204–9212, 2016.

- [20] J. Zhu, Y. Song, D. Jiang, and H. Song, "Multi-armed bandit channel access scheme with cognitive radio technology in wireless sensor networks for the Internet of Things," *IEEE Access*, vol. 4, pp. 4609–4617, 2016.
- [21] S. Maghsudi and S. Stańczak, "Hybrid centralized–distributed resource allocation for device-to-device communication underlying cellular networks," *IEEE Trans. Veh. Technol.*, vol. 65, no. 4, pp. 2481–2495, Apr. 2016.
- [22] J. Yang, B. Chen, J. Zhou, and Z. Lv, "A low-power and portable biomedical device for respiratory monitoring with a stable power source," *Sensors*, vol. 15, no. 8, pp. 19618–19632, 2015.
- [23] D. Jiang, Z. Xu, J. Liu, and W. Zhao, "An optimization-based robust routing algorithm to energy-efficient networks for cloud computing," *Telecommun. Syst.*, vol. 63, no. 1, pp. 89–98, Mar. 2015.
- [24] S. M. Sanchez, R. D. Souza, E. M. G. Fernandez, and V. A. Reguera, "Rate and energy efficient power control in a cognitive radio ad hoc network," *IEEE Signal Process. Lett.*, vol. 20, no. 5, pp. 451–454, May 2013.
- [25] G. Aniba and S. Aissa, "Cross-layer designed adaptive modulation algorithm with packet combining and truncated ARQ over MIMO Nakagami fading channels," *IEEE Trans. Wireless Commun.*, vol. 10, no. 4, pp. 1026–1031, Apr. 2011.
- [26] Q. Gao *et al.*, "Robust QoS-aware cross-layer design of adaptive modulation transmission on OFDM systems in high-speed railway," *IEEE Access*, vol. 4, pp. 7289–7300, 2016.
- [27] B. Ramachandran and S. Shanmugavel, "Mobility adaptive cross layer design for reliable route discovery in ad-hoc networks," in *Proc. Int. Conf. Wireless Commun. Sensor Netw.*, Allahabad, India, 2007, pp. 69–73.
- [28] J. Chen, T. Lv, and H. Zheng, "Joint cross-layer design for wireless QoS content delivery," *EURASIP J. Adv. Signal Process.*, vol. 2005, no. 2, pp. 1–16, 2005.
- [29] T. Praveena, G. S. Nagaraja, C. S. Shashank, and N. Reddy, "A novel scheduling algorithm emphasizing fairness for cross layer design in wireless networks," in *Proc. Int. Conf. Comput. Syst. Inf. Technol. Sustain. Solutions (CSITSS)*, Bengaluru, India, 2016, pp. 210–215.
- [30] Y. Dong and H. Dong, "Simulation study on cross-layer design for energy conservation in underwater acoustic networks," in *Proc. Oceans*, San Diego, CA, USA, 2013, pp. 1–5.
- [31] J. Zhu *et al.*, "Optimal and suboptimal access and transmission policies for dynamic spectrum access over fading channels in cognitive radio networks," *Chin. J. Electron.*, vol. 17, no. 4, pp. 726–732, 2008.
- [32] J. Zhu, Z. Peng, and F. Li, "A transmission and scheduling scheme based on W-learning algorithm in wireless networks," in *Proc. 8th Int. ICST Conf. IEEE Commun. Netw. China (CHINACOM)*, Guilin, China, 2013, pp. 85–90.
- [33] X.-H. Lin, Y. Tan, and J.-L. Zhang, "A MDP-based energy efficient policy for wireless transmission," *Syst. Eng. Electron.*, vol. 36, no. 7, pp. 1433–1438, 2014.
- [34] X. Zheng, Z. Cai, J. Li, and H. Gao, "A study on application-aware scheduling in wireless networks," *IEEE Trans. Mobile Comput.*, vol. 16, no. 7, pp. 1787–1801, Jul. 2017.
- [35] X. Zheng, Z. Cai, J. Li, and H. Gao, "Scheduling flows with multiple service frequency constraints," *IEEE Internet Things J.*, vol. 4, no. 2, pp. 496–504, Apr. 2017.
- [36] Z. Duan, W. Li, and Z. Cai, "Distributed auctions for task assignment and scheduling in mobile crowdsensing systems," in *Proc. 37th IEEE Int. Conf. Distrib. Comput. Syst. (ICDCS)*, Atlanta, GA, USA, 2017, pp. 635–644.
- [37] X. Zheng, Z. Cai, J. Li, and H. Gao, "An application-aware scheduling policy for real-time traffic," in *Proc. 35th IEEE Int. Conf. Distrib. Comput. Syst. (ICDCS)*, Columbus, OH, USA, 2015, pp. 421–430.
- [38] H. S. Wang and N. Moayeri, "Finite-state Markov channel—a useful model for radio communication channels," *IEEE Trans. Veh. Technol.*, vol. 44, no. 1, pp. 163–171, Feb. 1995.
- [39] V. K. N. Lau, "Performance of variable rate bit interleaved coding for high bandwidth efficiency," in *Proc. Veh. Technol. Conf.*, vol. 3, Tokyo, Japan, 2000, pp. 2054–2058.
- [40] S. T. Chung and A. J. Goldsmith, "Degrees of freedom in adaptive modulation: A unified view," *IEEE Trans. Commun.*, vol. 49, no. 9, pp. 1561–1571, Sep. 2001.
- [41] M. Humphrys, "W-learning: Competition among selfish Q-learners," *Comput. Lab., Univ. Cambridge, Cambridge, U.K., Tech. Rep. 362*, 1995.
- [42] Q. Wei, D. Liu, and G. Shi, "A novel dual iterative Q-learning method for optimal battery management in smart residential environments," *IEEE Trans. Ind. Electron.*, vol. 62, no. 4, pp. 2509–2518, Apr. 2015.
- [43] J. Ni, M. Liu, L. Ren, and S. X. Yang, "A multiagent Q-learning-based optimal allocation approach for urban water resource management system," *IEEE Trans. Autom. Sci. Eng.*, vol. 11, no. 1, pp. 204–214, Jan. 2014.
- [44] D. Silver *et al.*, "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [45] M. Spencer, J. Eickholt, and J. Cheng, "A deep learning network approach to ab initio protein secondary structure prediction," *IEEE/ACM Trans. Comput. Biol. Bioinformat.*, vol. 12, no. 1, pp. 103–112, Jan./Feb. 2015.
- [46] W. Huang, G. Song, H. Hong, and K. Xie, "Deep architecture for traffic flow prediction: Deep belief networks with multitask learning," *IEEE Trans. Intell. Transp. Syst.*, vol. 15, no. 5, pp. 2191–2201, Oct. 2014.
- [47] M. Wu and L. Chen, "Image recognition based on deep learning," in *Proc. Chin. Autom. Congr. (CAC)*, Wuhan, China, 2015, pp. 542–546.
- [48] W. Liu *et al.*, "Dispatching algorithm design for elevator group control system with Q-learning based on a recurrent neural network," in *Proc. Chin. Control Decis. Conf.*, Guiyang, China, 2013, pp. 3397–3402.
- [49] C. Li, J. Zhang, and Y. Li, "Application of artificial neural network based on Q-learning for mobile robot path planning," in *Proc. IEEE Int. Conf. Inf. Acquisition*, Weihai, China, 2006, pp. 978–982.
- [50] T. Kobayashi, T. Shibuya, F. Tanaka, and M. Morita, "Q-learning in continuous state-action space by using a selective desensitization neural network," *IEICE, Tokyo, Japan, Tech. Rep. 111*, 2011, pp. 119–123.
- [51] I. S. Comsa, S. Zhang, M. Aydin, P. Kuonen, and J.-F. Wagen, "A novel dynamic Q-learning-based scheduler technique for LTE-advanced technologies using neural networks," in *Proc. IEEE Conf. Local Comput. Netw.*, Clearwater, FL, USA, 2012, pp. 332–335.
- [52] T.-H. Teng and A.-H. Tan, "Fast reinforcement learning under uncertainties with self-organizing neural networks," in *Proc. IEEE/WIC/ACM Int. Conf. Web Intell. Agent Technol.*, Singapore, 2015, pp. 51–58.
- [53] C. B. Browne *et al.*, "A survey of Monte Carlo tree search methods," *IEEE Trans. Comput. Intell. Artif. Intell. Games*, vol. 4, no. 1, pp. 1–43, Mar. 2012.
- [54] C. J. C. H. Watkins and P. Dayan, "Technical note: Q-learning," *Mach. Learn.*, vol. 8, nos. 3–4, pp. 279–292, 1992.



Jiang Zhu received the Ph.D. degree in communication and information from the University of Electronic Science and Technology, Chengdu, China, in 2009.

He is an Associate Professor with the School of Telecommunications and Information Engineering, Chongqing University of Posts and Telecommunications, Chongqing, China. His current research interests include cognitive radio, wireless sensor networks, and mobile communication security technology.



Yonghui Song received the B.E. degree from the North China Institute of Aerospace Engineering, Langfang, China, in 2014. He is currently pursuing the M.E. degree at the Chongqing University of Posts and Telecommunications.

His current research interests include wireless sensor networks and cognitive radio.



Dingde Jiang (S'08–M'09) received the Ph.D. degree in communication and information systems from the University of Electronic Science and Technology of China, Chengdu, China, in 2009.

He is currently a Professor with the School of Astronautics and Aeronautic, University of Electronic Science and Technology of China, Chengdu. He was a Professor with the School of Computer Science and Engineering and the College of Information Science and Engineering, Northeastern University, Shenyang, China, prior to

2017. From 2013 to 2014, he was a Visiting Scholar with the Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN, USA. His research has been supported by the National Science Foundation of China, the Program for New Century Excellent Talents in University of Ministry of Education of China, etc. His current research interests include network measurement, modeling and optimization, performance analysis, network management, network security in communication networks, particularly in software-defined networks, information-centric networking, energy-efficient networks, and cognitive networks.

Dr. Jiang was a recipient of the Best Paper Award at several international conferences. He has been serving as an Editor for one international journal. He served as a Technical Program Committee member for several international conferences.



Houbing Song (M'12–SM'14) received the Ph.D. degree in electrical engineering from the University of Virginia, Charlottesville, VA, USA, in 2012.

In 2017, he joined the Department of Electrical, Computer, Software, and Systems Engineering, Embry-Riddle Aeronautical University, Daytona Beach, FL, USA, where he is currently an Assistant Professor and the Director of the Security and Optimization for Networked Globe Laboratory. He served on the faculty of West Virginia University, Morgantown, WV, USA, from 2012 to 2017. In

2007 he was an Engineering Research Associate with the Texas A&M Transportation Institute, Austin, TX, USA. He has edited four books including *Smart Cities: Foundations, Principles and Applications* (Wiley, 2017), *Security and Privacy in Cyber-Physical Systems: Foundations, Principles and Applications* (Wiley–IEEE Press, 2017), *Cyber-Physical Systems: Foundations, Principles and Applications* (Academic, 2016), and *Industrial Internet of Things: Cybermanufacturing Systems* (Springer, 2016). He has authored over 100 papers. His current research interests include cyber-physical systems, cybersecurity and privacy, Internet of Things, edge computing, big data analytics, unmanned aircraft systems, connected vehicle, smart and connected health, and wireless communications and networking.

Dr. Song was a recipient of the Golden Bear Scholar Award, the highest campus-wide recognition for research excellence at the West Virginia University Institute of Technology (WVU Tech), in 2016. He serves as an Associate Technical Editor for *IEEE Communications Magazine*. He is a Senior Member of the ACM.