

# JOINT TRANSACTION TRANSMISSION AND CHANNEL SELECTION IN COGNITIVE RADIO BASED BLOCKCHAIN NETWORKS: A DEEP REINFORCEMENT LEARNING APPROACH

Nguyen Cong Luong<sup>1</sup>, Tran The Anh<sup>1</sup>, Huynh Thi Thanh Binh<sup>2</sup>, Dusit Niyato<sup>1</sup>, Dong In Kim<sup>3</sup>, and Ying-Chang Liang<sup>4</sup>

<sup>1</sup>School of Computer Science and Engineering, Nanyang Technological University, Singapore

<sup>2</sup>School of Information and Communication Technology, Hanoi University of Science and Technology, Vietnam

<sup>3</sup>School of Information and Communication Engineering, Sungkyunkwan University, Korea

<sup>4</sup>Center for Intelligent Networking and Communications, University of Electronic Science and Technology of China, China

## ABSTRACT

To ensure that the data aggregation, data storage, and data processing are all performed in a decentralized but trusted manner, we propose to use the blockchain with the mining pool to support IoT services based on cognitive radio networks. As such, the secondary user can send its sensing data, i.e., transactions, to the mining pools. After being verified by miners, the transactions are added to the blocks. However, under the dynamics of the primary channel and the uncertainty of the mempool state of the mining pool, it is challenging for the secondary user to determine an optimal transaction transmission policy. In this paper, we propose to use the deep reinforcement learning algorithm to derive an optimal transaction transmission policy for the secondary user. Specifically, we adopt a Double Deep-Q Network (DDQN) that allows the secondary user to learn the optimal policy. The simulation results clearly show that the proposed deep reinforcement learning algorithm outperforms the conventional Q-learning scheme in terms of reward and learning speed.

**Index Terms**— Cognitive radio, blockchain, IoT, channel access, deep reinforcement learning

## 1. INTRODUCTION

Cognitive radio has been adopted to support IoT data transmission from IoT devices to a centralized server or the cloud [1]. Specifically, the IoT devices act as the Secondary Users (SUs) accessing idle spectrum of the Primary Users (PUs), improving the IoT performance and enhancing PUs' spectrum utilization. Generally, the IoT data to support IoT services and applications involves multiple stakeholders including devices owners, service providers, and users. Hence, the traditional approach of maintaining the IoT data by a single entity, e.g., an IoT provider, has shown many limitations. Firstly, it lacks transparency and traceability, i.e., the data can be modified arbitrarily by unknown persons and applications. Secondly, security is limited as it has to rely on a single entity which can be an easy target of cyber attacks. Thirdly, efficiency, speed, and reliability are low because of a bottleneck

and a single point of failure. This calls for a novel solution of the data management.

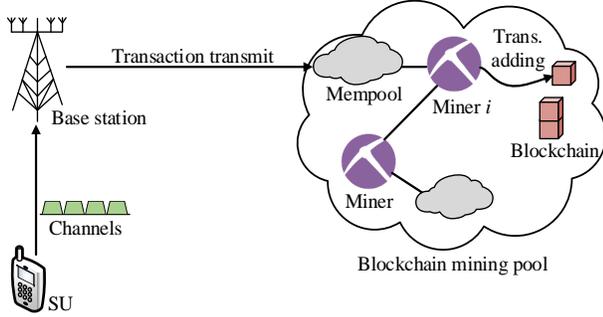
To overcome the limitations, we propose to use the blockchain [2] for collecting, storing, and processing the sensing data from the SUs. The first reason is that the blockchain is considered to be a decentralized database, i.e., a ledger [3] in which transactions, i.e., the sensing data, are recorded and processed by a number of nodes over the whole network instead of a centralized authority. The second reason is that the blockchain enhances the security and guarantees the data integrity since the transactions must be agreed and verified by the nodes before they are recorded [4]. Therefore, the blockchain can be combined with the cognitive radio to constitute a new framework called *cognitive radio based blockchain network*. The framework allows the SUs to use idle channels from the PUs to transmit their sensing data to the blockchain. The SU transmission is likely to be “localized” because of inherent Device-to-Device (D2D) transmission, which is well matched to P2P connection based blockchain network. Then, the SU's sensing data is recorded and processed in the blockchain in a decentralized but trusted manner.

However, under the dynamics of the primary channel and the uncertainty of the blockchain system, it may be challenging for each SU to make optimal decisions, i.e., transmit decision and channel selection, that maximizes the number of successful transaction transmissions. To address the challenge, we propose to use the Deep Q-Learning (DQL) technique presented in [5], i.e., the combination of Deep Neural Networks (DNNs) and the Q-Learning (QL) [6], that enables the SU to learn the optimal policy without requiring the prior information from the network environment. We first formulate an optimization problem for the SU that maximizes the number of successful transaction transmissions while minimizing the channel cost and transaction fee. Then, we adopt the Double Deep Q-Network (DDQN) to implement the DQL algorithm. Simulation results show that the proposed DQL outperforms the QL in terms of the performance and learning speed. To the best of our knowledge, this is the first work that studies the application of DQL [5] in the cognitive radio based blockchain

network.

The rest of this paper is organized as follows. Section 2 describes the system model. Section 3 presents the problem formulation. Section 4 presents the DQL algorithm for the joint transaction transmission and channel selection in the cognitive radio based blockchain network. Section 5 shows the performance evaluation results. Section 6 summarizes the paper.

## 2. SYSTEM MODEL



**Fig. 1.** A cognitive radio based blockchain network.

We consider a cognitive radio based blockchain network as shown in Fig. 1. The network consists of one SU, i.e., the IoT device, a base station, and the blockchain mining pool to support the IoT services. The base station can be regarded as a secondary receiver to establish a D2D connection with the SU. At each time slot  $t$ , the SU senses and selects one of  $K$  channels from multiple PUs to transmit a transaction to the mining pool through a base station. Here, the channel can be good, i.e., idle, or bad, i.e., busy, due to the transmission of the corresponding PU. The transaction contains sensing data of the SU and a transaction fee  $C_T$  that the SU is willing to pay the mining pool. After being successfully verified by miners in the mining pool, the transaction is stored in a mempool of one miner, say miner  $i$ , as an *unconfirmed transaction*. Note that the mempool is able to store  $D_{\max}$  unconfirmed transactions from the SUs in the network. The miner then adds the certain number of unconfirmed transactions with high transaction fees into a new block. Thus, the probability that the transaction of the SU is successfully added to the block at the current time slot is high if the transaction is assigned with a high fee [7]. The blockchain system is specifically vulnerable to the double-spending attack [8] in which transactions in a block can be maliciously modified. We assume that the transaction of the SU is attacked with a probability  $p_a$ .

## 3. PROBLEM FORMULATION

We formulate a stochastic optimization problem for the joint transaction transmission and channel selection of the SU. The problem is defined by a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R} \rangle$ , where  $\mathcal{S}$ ,  $\mathcal{A}$ , and  $\mathcal{R}$  are the state space, action space, and the reward function of the SU, respectively.  $\mathcal{P}$  is the state transition probabil-

ity function with  $P_{s,s'}(a)$  being the probability that the current state  $s \in \mathcal{S}$  transits to the next state  $s' \in \mathcal{S}$  when action  $a \in \mathcal{A}$  is executed.

### 3.1. Action Space

Let  $K$  denote the number of channels that the SU can choose to transmit its transactions. Then, the action space of the SU is defined as  $\mathcal{A} = \{0, 1, \dots, K\}$ , where  $a = 0$  means that the SU chooses not to transmit its transaction, and  $a = k$  means that the SU chooses channel  $k$  to transmit the transaction.

### 3.2. State Space

The state space is the combination of the channel state, denoted by  $\mathcal{S}^c$ , and the mempool state, denoted by  $\mathcal{S}^m$ .

First, we define  $\mathcal{S}^c$ . Each channel  $k$  can be in one of two different states, i.e., good or bad, i.e., the channel is idle or busy because of the transmission by the PU, respectively. The channels can be considered to be correlated, and thus all the channel states can be described as a  $2^K$ -state Markov chain [9] with a transition matrix  $\mathbf{P}$ . At the beginning of each time slot, although the SU cannot observe the states of all the channels, it can infer the states from its past channel selections, i.e., actions, and the corresponding observations.  $\mathcal{S}^c$  is thus defined as

$$\mathcal{S}^c = \{[a(t), w(t)], \dots, [a(t-L+1), w(t-L+1)]\}, \quad (1)$$

where  $w(t-l)$  is the observation of the channel selection at time slot  $t-l$ .  $w(t-l) = 1$  if the channel is good, and  $w(t-l) = 0$  if the channel is bad.  $L$  is the number of observations.

Second, we define mempool state  $\mathcal{S}^m$ .  $\mathcal{S}^m$  refers to the current number of transactions and the corresponding transaction fees in the mempool.  $\mathcal{S}^m$  is defined as  $\mathcal{S}^m = \{(m_1, \Delta C_1), \dots, (m_M, \Delta C_M)\}$ , where  $\Delta C_i$  represents transaction fee range  $i$ , and  $m_i$  is the current number of transactions that have transaction fees within the range of  $\Delta C_i$ . The transition of the mempool state from time slot  $t$  to  $t+1$  depends on (i) the number of transactions arriving in the mempool, (ii) the corresponding transaction fees, and (iii) the number of transactions that the miner adds a new block at time slot  $t$ . Here, we assume that the number of transactions arriving in the mempool and the transaction fee follow the uniform distributions  $U[T^{\min}, T^{\max}]$  and  $U[C_T^{\min}, C_T^{\max}]$ , respectively. The number of transactions added to the new block also follows the uniform distribution  $U[T_{\text{add}}^{\min}, T_{\text{add}}^{\max}]$ .

The state space of the SU at time slot  $t$  is thus defined as  $\mathcal{S} = \mathcal{S}^c \times \mathcal{S}^m$ , where  $\times$  is the Cartesian product.

### 3.3. Reward Function

The reward function  $\mathcal{R}$  of the SU is composed of three components, i.e., the positive utility  $R_{\text{success}}$ , the channel access

cost  $C_c$ , and the transaction fee  $C_T$ . The SU receives the utility of  $R_{\text{success}} > 0$  if the transaction transmission is successful and the utility of  $R_{\text{success}} = 0$  otherwise. The transaction transmission is considered to be “successful” if it is added to the new block at the current time slot and is not attacked. Here, we introduce the double-spending attack in which the transaction is attacked with a probability  $p_a$  given by [10]:

$$p_a = \begin{cases} 1 - \sum_{m=0}^n \binom{m+n-1}{m} (p^n q^m - p^m q^n) & \text{if } q < p, \\ 1 & \text{if } q \geq p, \end{cases}$$

where  $n$  and  $m$  are respectively the numbers of blocks that are found by the honest network and the attacker,  $p$  and  $q$ ,  $p + q = 1$ , are the probabilities that a block is found by the honest network and the attacker, respectively.

The objective is to maximize the number of successful transaction transmissions and minimize the channel cost and the transaction fee. Thus, the reward function of the SU is defined as  $\mathcal{R}(s, a) = R_{\text{success}} - C_c - C_T$ .

To obtain the mapping from a state  $s \in \mathcal{S}$  to an action  $a \in \mathcal{A}$  such that the long-term accumulated reward is maximized, the QL algorithm can be used. The algorithm finds the optimal policy defined as  $\pi^* : \mathcal{S} \rightarrow \mathcal{A}$  by estimating Q-values of state-action pairs, i.e.,  $Q(s, a)$ .  $Q(s, a)$  is the expected discounted sum of future rewards obtained by taking an action  $a$  at state  $s$  following the optimal policy. The Q-values are updated based on the experience of the SU as follows:

$$Q^{\text{new}}(s, a) = (1 - \lambda)Q(s, a) + \lambda \left( r(s, a) + \gamma \max_{a' \in \mathcal{A}} Q(s', a') \right), \quad (2)$$

where  $\lambda$  is the learning rate, and  $\gamma$  is the discount factor.

After the values  $Q(s, a)$  are learned, the SU can determine its optimal action from any state to maximize the long-term accumulated reward. However, the QL suffers from large state and action spaces of the network. Thus, we propose to use a DQL to find the optimal policy for the SU.

#### 4. DEEP Q-LEARNING ALGORITHM

DQL uses a DNN instead of the look-up table to represent all the states and actions of the SU. The input of the DNN is one of the states of the SU, and the output includes Q-values of all possible actions. To enable the SU to map its current state to an optimal action, the DNN needs to be trained. Training the DNN is to update its weights  $\theta$  by using experiences  $e = \langle s, a, r, s' \rangle$  of the SU to minimize a loss function. Here, the SU can execute action  $a$  using the  $\epsilon$ -greedy policy to balance its exploration and exploitation. The loss function at the current iteration is given by  $L = \mathbb{E} [(y(t) - Q(s, a; \theta))^2]$ , where  $y$  is the target value. Typically,  $y$  is defined as  $y = r + \gamma \max_{a' \in \mathcal{A}} Q(s', a'; \theta^{\text{old}})$ , where  $\theta^{\text{old}}$  are the weights of the DNN at the last iteration. However, such definition results

in over-optimistic value estimates since the max operator in  $y$  uses the same Q-values both to select and to evaluate an action. To decouple the action selection from the action evaluation, we propose to use the DDQN [11] which is composed of one online DNN with weights  $\theta^{\text{online}}$  and one target DNN with weights  $\theta^{\text{target}}$ . The online DNN updates its weights  $\theta^{\text{online}}$  at each iteration. The target DNN resets its weights  $\theta^{\text{target}}$  to  $\theta^{\text{online}}$  in every  $N^{\text{target}}$  iterations and keeps weights  $\theta^{\text{target}}$  fixed at other iterations. The online DNN updates its weights  $\theta$  to minimize the loss function defined as

$$L^{\text{DDQN}} = \mathbb{E} [(y^{\text{DDQN}} - Q(s, a; \theta^{\text{online}}))^2], \quad (3)$$

where the target value  $y^{\text{DDQN}}$  is defined as

$$y^{\text{DDQN}} = r + \gamma Q\left(s', \arg \max_{a' \in \mathcal{A}} Q_i(s', a'; \theta^{\text{online}}); \theta^{\text{target}}\right). \quad (4)$$

(4) shows that the selection of an action is due to the current weights, i.e.,  $\theta^{\text{online}}$ , while the weights  $\theta^{\text{target}}$  of the target DNN are used to evaluate fairly the value of the action.

---

**Algorithm 1** DQL algorithm [5].

---

**Input:**  $\mathcal{A}; N^{\text{target}}; N_b; \mathcal{M}$   
**Output:** Optimal policy  $\pi^*$

- 1: **Initialize:**  $\theta^{\text{online}}; \theta^{\text{target}}$
- 2: **for** episode  $i = \{1, \dots, N\}$  **do**
- 3:   **for** iteration  $t = \{1, \dots, T\}$  **do**
- 4:     Execute action  $a$  according to  $\epsilon - greedy$  policy
- 5:     Receive reward  $r_t$
- 6:     Store experience  $(s, a, r_t, s')$  in  $\mathcal{M}$
- 7:     Sample  $N_b$  experiences  $(s, a, r_j, s')$  from  $\mathcal{M}$
- 8:     **if** an episode terminates at iteration  $j + 1$  **then**
- 9:       Set  $y_j^{\text{DDQN}} = r_j$
- 10:    **else**
- 11:     Determine  $a^{\text{max}} = \arg \max_{a' \in \mathcal{A}} Q(s', a'; \theta^{\text{online}})$
- 12:     Set  $y_j^{\text{DDQN}} = r_j + \gamma Q(s', a^{\text{max}}; \theta^{\text{target}})$
- 13:    **end if**
- 14:    Perform a gradient descent step on  $(y_j^{\text{DDQN}} - Q(s, a; \theta^{\text{online}}))^2$  to update  $\theta^{\text{online}}$
- 15:    Reset  $\theta^{\text{target}} = \theta^{\text{online}}$  in every  $N^{\text{target}}$  iterations
- 16:    **end for**
- 17: **end for**

---

Algorithm 1 shows the DQL algorithm which uses the DDQN to find the optimal policy for the SU. Accordingly, based on the experience  $e$ , the online DNN and target DNN compute the optimal value  $Q(s', a'; \theta^{\text{online}})$ . Then, the target value  $y^{\text{DDQN}}$  and the loss function  $L^{\text{DDQN}}$  is calculated according to (4) and (3), respectively. The value of  $L^{\text{DDQN}}$  is used to update weights  $\theta$  of the online DNN. To ensure the stability of the learning, the experience replay memory  $\mathcal{M}$  is used to store experience  $e$ , and then a mini-batch of  $N_b$  experiences are taken at each iteration to train the DNNs.

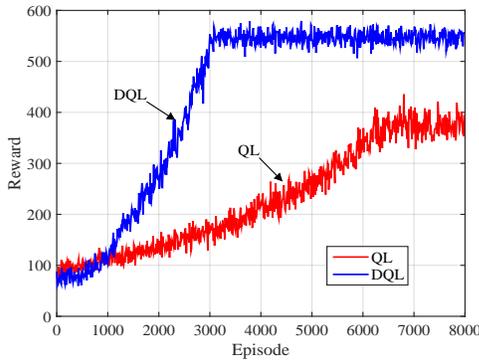
#### 5. PERFORMANCE EVALUATION

In this section, we present experimental results to evaluate the performance of the proposed DQL algorithm. For compari-

son, the QL algorithm [6] is used as a baseline scheme. Major simulation parameters are listed in Table 1. The simulation results for the performance comparison between the proposed DQL scheme and the QL scheme are shown in Figs. 2, 3, 4, and 5 depending on the varied parameters.

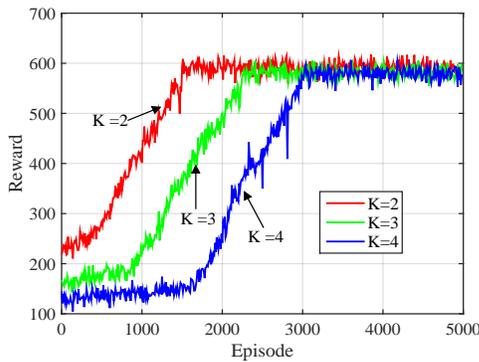
**Table 1.** Simulation parameters

Parameters	Value
Number of channels ( $K$ )	4
Probability of switching good channel ( $p_c$ )	0.9
Maximum number of transactions in the mempool ( $D_{\max}$ )	50
Channel cost ( $C_c$ )	0.2
Transaction fee ( $C_T$ )	$\sim U[0; 1]$
Probability that a block is found by the attacker ( $q$ )	0.02
Discount rate ( $\gamma$ )	0.9
$\epsilon$ -greedy	$0.9 \rightarrow 0$



**Fig. 2.** DQL scheme and QL scheme comparison.

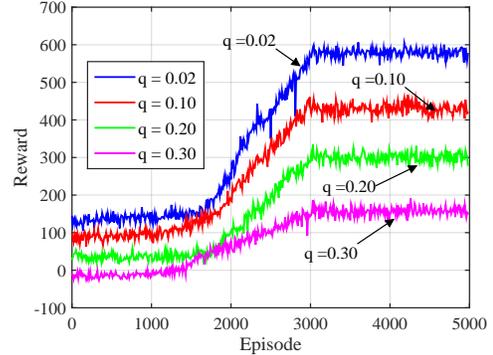
Fig. 2 illustrates the rewards obtained by the DQL and QL schemes. To enable the QL scheme to be run in our computation environment, we reduce the state space by setting the maximum number of transactions in the mempool to be 10. As seen, the DQL scheme converges to the reward much higher than that of the QL scheme. Specifically, the rewards obtained by the DQL and QL schemes are 550 and 390, respectively. Moreover, the convergence speed of the DQL scheme is faster than that of the QL scheme. The DQL scheme converges at around 3000 episodes while the QL scheme converges at 7000 episodes.



**Fig. 3.** Reward of DQL as the number of channels  $K$  is varied.

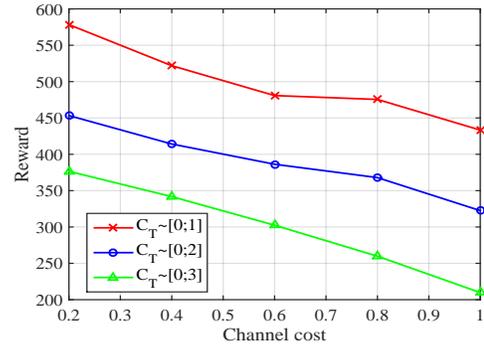
The convergence speed of the DQL scheme is likely main-

tained as the maximum number of transactions in the mempool increases to 50 as shown in Fig. 3. In this case, the state space becomes too large for the QL scheme to converge in the reasonable time, and hence it is not shown in the figure. This confirms the scalability of the DQL. Note that as the number of channels  $K$  is varied, the state space changes, and the DQL scheme has different convergence speeds. However, the DQL scheme always reaches to the same reward because it already learns the optimal policy to obtain the maximum reward.



**Fig. 4.** Reward of DQL as the probability  $q$  is varied.

The reward obtained by the DQL scheme decreases as the probability that a block is found by the attacker  $q$  increases as shown in Fig. 4. The reason is that as  $q$  increases, the number of successfully transmitted transactions decreases. Similarly, the reward that the SU receives decreases as the transaction fee  $C_T$  and the channel cost  $C_c$  increase as illustrated in Fig. 5.



**Fig. 5.** Reward of DQL versus the channel access cost  $C_c$ .

## 6. CONCLUSIONS

In this paper, we have presented the DQL algorithm for the joint transaction transmission and channel selection problem in the cognitive radio based blockchain network. Specifically, we have developed a DQL algorithm using DDQN to solve the problem. The simulation results show that the proposed DQL scheme outperforms the QL scheme in terms of reward and learning speed. This implies that the DQL enables the SU to achieve the higher number of successful transaction transmissions while paying lower cost.

## 7. REFERENCES

- [1] A. A. Khan, M. H. Rehmani, and A. Rachedi, "Cognitive-radio-based internet of things: Applications, architectures, spectrum related functionalities, and future research directions," *IEEE wireless communications*, vol. 24, no. 3, pp. 17–25, June 2017.
- [2] O. Schrijvers, J. Bonneau, D. Boneh, and T. Roughgarden, "Incentive compatibility of bitcoin mining pool reward functions," in *International Conference on Financial Cryptography and Data Security*. Barbados: Springer, May 2016, pp. 477–498.
- [3] R. Neisse, G. Steri, and I. Nai-Fovino, "A blockchain-based approach for data accountability and provenance tracking," in *Proceedings of the 12th International Conference on Availability, Reliability and Security*. Reggio Calabria, Italy: ACM, August 2017, pp. 93–98.
- [4] B. Liu, X. L. Yu, S. Chen, X. Xu, and L. Zhu, "Blockchain based data integrity service framework for iot data," in *IEEE International Conference on Web Services (ICWS)*, Honolulu, HI, June 2017, pp. 468–475.
- [5] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [6] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, no. 3-4, pp. 279–292, May 1992.
- [7] (2018, Sept.) Explaining bitcoin transaction fees. [Online]. Available: <https://support.blockchain.com/hc/en-us/articles/>
- [8] G. O. Karame, E. Androulaki, and S. Capkun, "Double-spending fast payments in bitcoin," in *ACM conference on Computer and communications security*, Raleigh, NC, October 2012, pp. 906–917.
- [9] S. Wang, H. Liu, P. H. Gomes, and B. Krishnamachari, "Deep reinforcement learning for dynamic multichannel access in wireless networks," *IEEE Transactions on Cognitive Communications and Networking*, vol. 4, no. 2, pp. 257–265, Jun. 2018.
- [10] M. Rosenfeld, "Analysis of hashrate-based double spending," *arXiv preprint arXiv:1402.2009*, 2014.
- [11] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*. Phoenix, AZ, February 2016, pp. 2094–2100.