

Digital Twin-Assisted Efficient Reinforcement Learning for Edge Task Scheduling

Xiucheng Wang*, Longfei Ma*, Haocheng Li*, Zhisheng Yin†, Tom. Luan†, and Nan Cheng*

*School of Telecommunications Engineering, Xidian University, Xi'an, China

†School of Cyber Engineering, Xidian University, Xi'an, China

Email: {xcwang_1, lfma, lhc}@stu.xidian.edu.cn, {zsyin, tom.luan}@xidian.edu.cn
, dr.nan.cheng@ieee.org

Abstract—Task scheduling is a critical problem when one user offloads multiple different tasks to the edge server. When a user has multiple tasks to offload and only one task can be transmitted to server at a time, while server processes tasks according to the transmission order, the problem is NP-hard. However, it is difficult for traditional optimization methods to quickly obtain the optimal solution, while approaches based on reinforcement learning face with the challenge of excessively large action space and slow convergence. In this paper, we propose a Digital Twin (DT)-assisted RL-based task scheduling method in order to improve the performance and convergence of the RL. We use DT to simulate the results of different decisions made by the agent, so that one agent can try multiple actions at a time, or, similarly, multiple agents can interact with environment in parallel in DT. In this way, the exploration efficiency of RL can be significantly improved via DT, and thus RL can converge faster and local optimality is less likely to happen. Particularly, two algorithms are designed to make task scheduling decisions, i.e., DT-assisted asynchronous Q-learning (DTAQL) and DT-assisted exploring Q-learning (DTEQL). Simulation results show that both algorithms significantly improve the convergence speed of Q-learning by increasing the exploration efficiency.

Index Terms—task scheduling, digital twin, reinforcement learning, exploration efficiency

I. INTRODUCTION

With the development of information technologies such as Internet technology, artificial intelligence (AI), and computer vision, users are increasingly demanding services such as image recognition or VR/AR that require strong computing capability. Since such services usually have a large task size, offloading all these tasks to the cloud will result in an overburdened communication. Therefore, offloading tasks to the server in edge is emerged as an promising solution to this issue. On the other hand, due to the deployment on edge nodes, the computing capability of edge servers is limited. As a result, edge servers usually processes only a few (very likely one) tasks at the same time to provide ensuring service quality. Meanwhile, user has multiple tasks with different delay requirements that need to be offloaded, and only one task can be transmitted at a time, the order of task offloading will directly affect the quality of service. In [1] scheduling tasks is proved a flow shop scheduling problem, and in [2] this kind of problem is proved as NP-Hard. Thus, it is difficult to optimize this problem with high accuracy and speed by using tradition optimization method.

To get optimal or near-optimal solutions quickly, RL-based on task sorting algorithms have attracted increasing attention. Monte carlo tree search is applied in the specific context of task scheduling [3]. However, there is a very large action space in task scheduling problem, e.g., making it difficult for RL to achieve excellent performance. This is because usual RL agent uses *exploration – exploitation* to converge. When agent use *exploration*, it chooses an random action through a specific distribution to obtain the performance of different permutations, and when agent use *exploitation*, it samples the best performance action repeatedly to make the algorithm converge. Obviously, *exploration* and *exploitation* are contradictory, the higher exploration probability, the slower it is to converge, while low exploration probability causes agent has no knowledge of enough actions, leading to a poor converge performance. Therefore, we need a better *exploration* and *exploitation* that has efficient exploration and high convergence speed.

DT is the virtual realization of entities in the real world through digital form. Unlike traditional simulation models, digital twins are the dynamic panoramic mapping of physical entities. Through data interaction with the entity, it constantly updates its own information to ensure the real-time, accurate and comprehensive virtual mapping, so as to realize the clone with high approximation on the virtual platform [4]. Furthermore, the virtual model can integrate various data from the physical world, and use expert knowledge, AI and other means to conduct comprehensive analysis, so as to better predict future states. Therefore, the virtual body in the digital twin is not only a simple copy of the information of the physical entity, but a complete clone of the state, characteristics and development trend of the physical entity. When a clone is established, various complex actions can be quickly performed on the virtual platform, and obtain the state of the system after taking these actions, so as to provide support for decision-making [5]. Compared with common methods of operating in the real world, digital twins can significantly improve efficiency and reduce costs.

Since DT can reproduce all the properties of real physical space in digital space, it enables one agent to try different actions at a time in a digital space, and also enables multiple agents to interact with the same environment in parallel come true. Therefore, to improve RL exploration efficiency, we

introduce DT to get a higher performance RL with lower time. In this paper, we use DT to enable the agent to learn performance of multiple actions at a time, thus increasing the exploiting efficiency of RL-based task scheduling. The proposed methods show better performance and faster convergence speed than traditional Q-learning based RL method. To our best knowledge, this is the first paper using DT to enhance RL performance, which not only gives a way to further optimize RL-based decision methods, but also shows a potential research direction for DT.

II. SYSTEM MODEL

We consider a single-user edge computing communication system [1], where the server is deployed near the user to provide service. At the beginning of each frame, the user generates N independent tasks and the task i is denoted by $t_i, i \in 1, 2, \dots, N$. Due to the limited computing power and energy of user, all tasks need to be transmitted to the edge server for execution. In this paper, we assume that the server uses a single-core CPU with constant CPU frequency, which can only process one task at a time. After the server completes a task, it immediately transmits the result back to the user.

A. Task and Communication Model

We denote the set of tasks of each communication frame by $\Gamma = \{t_1, t_2, \dots, t_N\}$. Each task is described by a ternary vector $\mathbf{t}_i = [d_i, \varepsilon_i, c_i]$, where d_i denotes the amount of the task data, ε_i is the task deadline, and c_i represents the task complexity which means the amount of CPU operations required to process 1 bit of data. In particular, ε_i represents the maximum delay can be tolerated from task generation to result reception, which reflects the urgency of the task. The user expects to receive all the results within the deadline, otherwise the timeout tasks will be regarded as failure.

In order to avoid useless waiting time, the server executes the tasks following the order of arrival. Assuming user can only transmit the data for one task at each time. As there is only one user in this edge computing system, we consider that the channel state, transmission power and distance are known at the service, so that the transmission rates of different tasks are the same, which is denoted by R .

In the case of busy CPU, the task arriving at the sever will wait in the memory queue. The memory of the server only stores the data of tasks in the current communication frame. Since the amount of tasks generated by a single user in a frame is limited, we consider the memory of the server is large enough so that there will be no data overflow. Therefore, the order of task execution is the same as the order of data transmission. The computing frequency of CPU is denoted by f_{ser} . In addition, since the size of result data is much smaller than the input data, we ignore the result transmission delay. For the user, the most important concern is the task completion rate and total completion delay, which depends on the queue order. We denote the queue of N tasks as $\sigma = [\sigma_1, \sigma_2, \dots, \sigma_N]$, where σ_i is the task of order i , $\sigma_i \in \{1, \dots, N\}$. Consequently, our optimization goal is to ensure that all tasks can be completed

within the deadline, and to minimize the total completion delay.

B. Problem Formulation

The completion time of task t_j is denoted by $T_{\text{comp}}^j(\sigma)$, which includes the delay of transmission, execution, and queuing. The execution of a task only begins when all of its data has been received by the server and there are no other tasks ahead of the queue. So only when t_j is at the head, it will be transmitted and executed without any wait. Other tasks must wait before transmission. We denote the whole data input time of the j -th task by $T_{\text{ready}}^j(\sigma)$, which is given by

$$T_{\text{ready}}^j(\sigma) = \sum_{i=1}^j \frac{d_{\sigma_i}}{R}, j = 1, \dots, N. \quad (1)$$

Obviously, $T_{\text{comp}}^j(\sigma)$ includes $T_{\text{ready}}^j(\sigma)$ and the delay in the server which depends on $T_{\text{ready}}^j(\sigma)$ and the completion time of the task ahead. Thus, it can be determined by

$$T_{\text{comp}}^j(\sigma) = \begin{cases} T_{\text{ready}}^j(\sigma) + d_{\sigma_j} c_{\sigma_j} f_{\text{ser}}^{-1}, & j = 1 \\ \max \left\{ T_{\text{ready}}^j(\sigma), T_{\text{comp}}^{j-1}(\sigma) \right\} \\ \quad + d_{\sigma_j} c_{\sigma_j} f_{\text{ser}}^{-1}, & j > 1 \end{cases} \quad (2)$$

where d_{σ_j} and c_{σ_j} represent the data size and complexity of t_{σ_j} respectively, and when $j = N$, $T_{\text{comp}}^j(\sigma)$ is the total completion time of the task sequence.

In order to simultaneously optimize the task success rate and total completion time within a frame, we formulate the queuing problem as

$$\min_{\sigma} \sum_{i=1}^N T_{\text{comp}}^i + \zeta \mathbb{I}(T_{\text{comp}}^i > \varepsilon_i), \quad (3)$$

where $\mathbb{I}(\ast)$ is indicator function, when \ast is true $\mathbb{I}(\ast) = 1$, otherwise $\mathbb{I}(\ast) = 0$, ζ is importance factor.

III. DT-ASSISTED Q-LEARNING METHOD

A. Q-learning

Reinforcement learning has become a powerful means to solve the problem of resource management and task scheduling in wireless communication networks [6]–[8], where Q-learning [9] is a value-based algorithm which has good convergence. The goal of reinforcement learning is to get the optimal policy which means helping the agent maximize the reward value. In most cases, the agent needs to take a series of actions to complete a task and the reward are delayed, so the expected total reward Q is seen as the evaluation function of the action, which is given by $Q = \sum_{t=1}^{\infty} \gamma^{t-1} r_t(s_t, a_t)$, where γ is the discount factor, which represents the influence of future states on the current policy, s_t and a_t are the state and action at time t , respectively, $r_t(s_t, a_t)$ is the present reward for taking a_t in s_t . Q-learning stores Q values by building a Q-Table whose coordinates are states and actions. With the table, the agent only needs to select the action with largest Q value

to execute according to the state of environment. However, it is often difficult to obtain the accurate Q-Table due to the unknowns of the environment.

For large state spaces, Q-learning uses Temporal-Difference method to approximate the true Q value. Specifically, when the agent observes state s_t at time t , the action is selected by the ϵ -greedy strategy, where in the greedy decision, the optimal action is given by

$$a_t^{opt} = \arg \max_{a_i \in A_t} Q(s_t, a_i), \quad (4)$$

where a_t^{opt} and A_t are the optimal action and the set of actions at time t , respectively. Execute the selected action and enter the next state, then update the Q value according to the feedback of the environment, which is formulated as

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1}(s_t, a_t) + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)], \quad (5)$$

where α is the learning rate, which is used to balance the efficiency and stability of learning. Based on the above method, iterate the Q value until get a reliable Q-Table.

A key point of Q-learning is that the agent adopts the ϵ -greedy strategy to choose actions, which can balance *exploitation* and *exploration* well. In details, ϵ -greedy means that when the agent makes a decision, there is a small probability ϵ ($\epsilon < 1$) to randomly select an action, and choose the largest-value action with the probability of remaining $1 - \epsilon$. Assume that the initial state is s_1 , the set of available actions is A_1 and the known optimal action is a_1^{opt} . After the agent takes an action, it receives a reward r_1 from environment. In the decision-making process, the probability of each non-optimal action being selected is $\frac{\epsilon}{|A_1|}$, where $|A_1|$ represents the number of actions. The probability of choosing a_1^{opt} is $\frac{\epsilon}{|A_1|} + 1 - \epsilon$. Therefore, by using the ϵ -greedy strategy, the agent can make a good trade-off between *exploitation* and *exploration*.

B. DT-Assisted Method

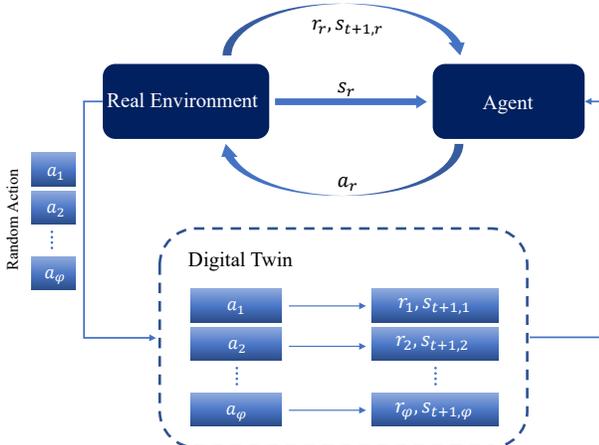


Fig. 1. DT-assisted exploring Q-learning algorithm flowchart.

Algorithm 1 Digital Twin-Assisted Asynchronous Q-learning Algorithm

```

1: Initialize all elements in Q-Table as 0, learning rate  $lr$ ,
   and update date cycle  $\delta$ 
2: for  $i = 1$  to epoch do
3:    $\epsilon = \epsilon_{min} + \epsilon e^{-i\beta}$ 
4:   With probability  $\epsilon$  select a random action  $a_r$ , otherwise
   select  $a_r = \arg \max_a Q(s_t, a)$ 
5:   Take action  $a_r$  in real environment and get reward  $r_r$ 
   and next state  $s_{t+1,r}$ 
6:    $Q_r(s_t, a_r) = Q(s_t, a_r) + lr(r + \gamma \max_a Q_r(s_{t+1}, a) -$ 
    $Q_r(s_t, a_r))$ 
7:   for  $j = 1$  to  $\phi$  do
8:      $j$ -th agent select a random action  $a_j$  with probability
      $\epsilon$ , otherwise select  $a = \arg \max_a Q_j(s_t, a)$ 
9:      $Q_j(s_t, a_j) = Q_j(s_t, a_j) + lr(r +$ 
      $\gamma \max_a Q_j(s_{t+1}, a) - Q_j(s_t, a_j))$ 
10:    end for
11:    if  $i \% \delta == 0$  then
12:       $Q = \frac{Q + \sum_{j=0}^{\phi} Q_j}{1 + \phi}$ 
13:       $Q_r = Q$ 
14:      for  $j = 1$  to  $\phi$  do
15:         $Q_j = Q$ 
16:      end for
17:    end if
18:  end for

```

Algorithm 2 Digital Twin-Assisted Exploring Q-learning Algorithm

```

1: Initialize  $Q(s, a)$  as 0 and learning rate  $lr$ 
2: for  $i = 1$  to epoch do
3:    $\epsilon = \epsilon_{min} + \epsilon e^{-i\beta}$ 
4:   With probability  $\epsilon$  select a random action  $a_r$ 
5:   otherwise select  $a_r = \arg \max_a Q(s_t, a)$ 
6:   Take action  $a_r$  in real environment and get reward  $r_r$ 
   and next state  $s_{t+1,r}$ 
7:   Randomly select  $\phi$  unique actions  $a_1, a_2, \dots, a_\phi$ 
8:   Take these action respectively in DT to get reward
    $r_1, r_2, \dots, r_\phi$  and next state  $s_{t+1,1}, s_{t+1,2}, \dots, s_{t+1,\phi}$ 
9:
10:  /* Update Q-Table */
11:   $Q(s_t, a_r) = Q(s_t, a_r) + lr(r + \gamma \max_a Q(s_{t+1}, a) -$ 
    $Q(s_t, a_r))$ 
12:  for  $j = 1$  to  $\phi$  do
13:     $Q(s_t, a_j) = Q(s_t, a_j) + lr(r + \gamma \max_a Q(s_{t+1}, a) -$ 
     $Q(s_t, a_j))$ 
14:  end for
15: end for

```

Traditional Q-learning can only interact with the real physical environment and can only explore one action at a time, so the exploration efficiency is low. In addition, since the real physical environment can only accept decisions made by one agent, the agent can only converge in one direction, and it

is easy to fall into a local optimum. Inspired by [10] we first tried to place multiple agents in DT, which means DT-assisted asynchronous Q-learning (DTAQL). Assuming that the simulation capability of DT is ϕ , which means DT can complete the simulation of ϕ times interaction between the agent and the environment within a tolerable time. Therefore, we assume there are ϕ agents in DT, and each agent maintains a Q-Table, and independently selects actions in DT to interact with the environment according to the ϵ -greedy algorithm. Then, each agent updates its Q-Table according to the actions it takes and the rewards the get. When all agents have updated the Q-Table δ times, they share their knowledge and update the Q-Table. Since all agents are independent of each other, they may take different actions, which improves the efficiency of exploration. Besides, through periodically sharing knowledge, the agents in the real physical environment can obtain rewards of different scheduling order faster, thereby improving the convergence speed. Because knowledge sharing is periodic, each agent may take completely different actions in one cycle and update its own Q-Table independently. This means that in one cycle, the convergence direction of different agents may be different, which is beneficial for the agents in the real physical environment to avoid getting stuck in local optima and obtain better convergence performance.

However, asynchronous reinforcement learning needs to maintain multiple Q-Tables, which is memory consuming. Due to periodically sharing knowledge, agents are not completely independent. This leads to the fact that although the convergence directions of the agents in one cycle may be different, but on a large scale all agents still converge in the same direction. Thus, the efficiency of exploration is limited, and the convergence speed and performance is degraded.

Therefore, we proposed another DT-assisted exploring Q-learning (DTEQL) method. In this method, only one agent is needed. First, the agent selects an action a_r to act on the real environment according to the ϵ -greedy algorithm, and obtains the r_r and $s_{t+1,r}$ of the real environment feedback. At the same time, randomly select different ϕ actions a_1, a_2, \dots, a_ϕ , and apply these ϕ actions to the virtual environment of DT respectively to obtain reward r_1, r_2, \dots, r_ϕ and next state $s_{t+1,1}, s_{t+1,2}, \dots, s_{t+1,\phi}$. Then, according to Equation 5, updates the Q value of these actions. Compared with DT-assisted asynchronous Q-learning, this method is more random in action selection thus improving the exploration efficiency. Moreover, because the agent always interacts with the real environment according to the ϵ -greedy algorithm, the probability of selecting the optimal action in the Q-Table is increased, so that the agent can converge to the optimal value faster.

IV. SIMULATION RESULT

In this section, we provide the simulation result to show the efficiency for DT-assisted reinforcement learning method. In simulation, task size d , task complexity c , and deadline ε are all assumed to distribute uniformly, $d \sim Unif[0, 2Mb]$, $c \sim Unif[0, 1000]$ CPU cycles and $\varepsilon \sim Unif[1, 5]s$. The

CPU frequency of edge server is 10GHz. We randomly generated 1,000 tasks to evaluate the convergence speed and performance for different algorithms. The learning rate lr is 0.1. We use the following algorithm for comparison.

- **QL**: The traditional Q-learning method using ϵ -greedy to explore the environment. There is only one agent in the real physical space, and the agent can only choose one action to interact with the environment at a time. Due to the huge state space, we set the minimum exploration rate ϵ_{min} as 0.1 and the exploration decay factor β as 5,000.

- **DTAQL**: DT-assisted asynchronous Q-learning method. Besides real physical space, ϕ agent is placed in the digital space to interact with the virtual environment and periodically share knowledge, and more details are presented in Algorithm 1. But when testing performance, we only test the performance of the agent in the real environment. As for δ , all agent share their knowledge every 512 iterations.

- **DTEQL**: DT randomly selects p different actions and simulates the outcomes of these actions, thereby increasing the exploration probability, and more details are presented in Algorithm 2. Similar to DTAQL, we only test the performance of the agent in the real environment.

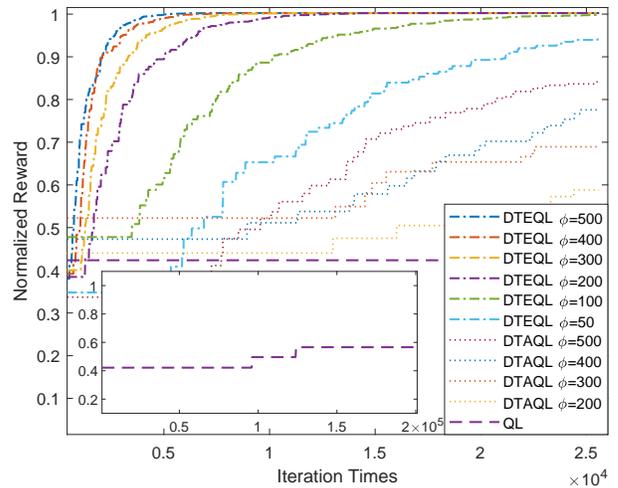


Fig. 2. Convergence performance of different algorithm with $N = 8$. Since QL need over 100,000 iterations to converge, while other algorithms need far much fewer iterations, so we show it in a separate figure specially.

TABLE I
PERFORMANCE OF THE DIFFERENT ALGORITHM

	DTEQL $\phi=500$	DTEQL $\phi=200$	DTEQL $\phi=100$	DTEQL $\phi=50$	DTAQL $\phi=500$	QL
Normalized Reward	1.0024	1.0024	0.9979	0.9406	0.7758	0.4235
Deadline Miss Ratio	0.0195	0.0195	0.0195	0.0273	0.0429	0.039
Average Delay	0.6435	0.6435	0.6453	0.6677	0.7321	0.8697
Convergence Time	5597	12752	25549	>25600	>25600	≥25600

First, we evaluate the performance of DTAQL when $N = 8$, under different number of agents ϕ in the DT. Obvious, as is shown in Fig 2, the larger ϕ is, the better the performance and fast convergence speed can be achieved. This is because different agents will select actions in parallel when exploring the environment [10]. As the number of agents ϕ in DT increases, the possibility of agents exploring more actions

TABLE II
TIME TO CONVERGE FOR THE DIFFERENT ALGORITHMS

	DTEQL $\phi=500$	DTEQL $\phi=200$	DTEQL $\phi=100$	DTEQL $\phi=50$	DTAQL $\phi=500$	QL
$N=6$	102	295	510	1047	1284	>25600
$N=7$	755	2497	3873	8795	10239	>>25600
$N=8$	5597	12752	25549	>25600	>25600	>>25600

will also increase, so that through periodic knowledge sharing improves the decision-making ability of the agent in the real physical environment. However, due to the periodic sharing of knowledge between agents, each agent actually has a similar exploration direction on a large scale, so the increase of agents does not significantly improve performance. Moreover, because knowledge sharing is periodic, the performance of agents in real physical environments shows a step-up trend.

Then, we test the performance of DTEQL. Similar to DTAQL, the convergence speed increases as ϕ increases. But when ϕ is with the high value region, as ϕ increases, the rate at which the convergence speed improves gradually decreases. As is shown in Fig 2, DTEQL converges faster for DTAQL even when ϕ is small than DTAQL with $\phi = 500$. The reason is that DTEQL randomly select ϕ different actions, thereby improving the exploration efficiency, and since the agent in the real physical environment still likes the ordinary Q-learning method, using the ϵ -greedy to select optimal action, which increases the sampling probability of the optimal action and improves the convergence speed. Besides, because DTEQL does not need to store multiple Q tables, its storage consumption is also less than DTAQL.

In Table I, we show specific data on reward, average task complete delay, deadline miss ratio, and number of times required to converge for different algorithms, if no improvement in algorithm performance can be observed during training, we use $\gg 25600$ to denote it. Limited by training time, we only trained 25,600 times for each algorithm, which is also a very long time. For DTEQL, when ϕ is greater than 200, the algorithm can always converge within 25,600 training times, and the performance after convergence is exactly the same. But as ϕ increases, the convergence speed also increases. When ϕ is 100, we can see that the reward has decreased, but deadline miss ratio keeps constant. This shows that although the task processing delay increases when DTEQL fails to converge completely, the algorithm can still ensure that all tasks are completed within the deadline. Although the performance of DTAQL is obviously better than that of the ordinary Q-learning method, it is much weaker than DTEQL in this problem.

Table II shows convergence time for different algorithm under different N . As the task number N decreases, the action space of the scheduling problem decreases, and the algorithm requires smaller ϕ to converge. This means we can adaptively change ϕ of DT, so as to reduce the construction cost of DT, since stronger DT means more accurate data and more advanced building technology [11], which inevitably bring higher price.

V. CONCLUSION

In this paper, we have investigated the task scheduling for edge offloading with the assistance of DT. By using DT to enrich the action space, we have proposed two DT-assisted RL algorithms to let the agent try many actions at the same time or multiple agents independently interact with the environment and exchange their knowledge periodically. Simulation results have shown DT can significantly assist improving the exploration efficiency, thereby the convergence speed of Q-learning can be increased and its convergence performance can be improved. Besides, users can experience a higher quality of service with lower latency by our proposed scheme. For future research, we will study the performance of DT-assisted task scheduling in more complex network such as space-air-ground integrated networks, whose action space is continuation.

ACKNOWLEDGEMENT

This work was supported by the National Key Research and Development Program of China (2020YFB1807700), the National Natural Science Foundation of China (NSFC) under Grant No. 62071356, and the Fundamental Research Funds for the Central Universities under Grant No. JB210113.

REFERENCES

- [1] Y. Mao, J. Zhang, and K. B. Letaief, "Joint task offloading scheduling and transmit power allocation for mobile-edge computing systems," in *2017 IEEE wireless communications and networking conference (WCNC)*. IEEE, 2017, pp. 1–6.
- [2] L. A. Herrbach and J. Y.-T. Leung, "Preemptive scheduling of equal length jobs on two machines to minimize mean flow time," *Operations Research*, vol. 38, no. 3, pp. 487–494, 1990.
- [3] Z. Hu, J. Tu, and B. Li, "Spear: Optimized dependency-aware task scheduling with deep reinforcement learning," in *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2019, pp. 2037–2046.
- [4] F. Tao, H. Zhang, A. Liu, and A. Y. C. Nee, "Digital twin in industry: State-of-the-art," *IEEE Transactions on Industrial Informatics*, vol. 15, no. 4, pp. 2405–2415, 2019.
- [5] W. Sun, H. Zhang, R. Wang, and Y. Zhang, "Reducing offloading latency for digital twin edge networks in 6g," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 10, pp. 12 240–12 251, 2020.
- [6] N. Cheng, F. Lyu, W. Quan, C. Zhou, H. He, W. Shi, and X. Shen, "Space/aerial-assisted computing offloading for iot applications: A learning-based approach," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 5, pp. 1117–1129, 2019.
- [7] C. Zhou, W. Wu, H. He, P. Yang, F. Lyu, N. Cheng, and X. Shen, "Deep reinforcement learning for delay-oriented iot task scheduling in sagin," *IEEE Transactions on Wireless Communications*, vol. 20, no. 2, pp. 911–925, 2021.
- [8] H. Huang, Q. Ye, and Y. Zhou, "Deadline-aware task offloading with partially-observable deep reinforcement learning for multi-access edge computing," *IEEE Transactions on Network Science and Engineering*, pp. 1–1, 2021.
- [9] B. Jang, M. Kim, G. Harerimana, and J. W. Kim, "Q-learning algorithms: A comprehensive classification and applications," *IEEE Access*, vol. 7, pp. 133 653–133 667, 2019.
- [10] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *International conference on machine learning*. PMLR, 2016, pp. 1928–1937.
- [11] M. Liu, S. Fang, H. Dong, and C. Xu, "Review of digital twin about concepts, technologies, and industrial applications," *Journal of Manufacturing Systems*, vol. 58, pp. 346–361, 2021.