

# Deep Reinforcement Learning for Backscatter-Aided Data Offloading in Mobile Edge Computing

Shimin Gong, Yutong Xie, Jing Xu, Dusit Niyato, and Ying-Chang Liang

## Abstract

Wireless network optimization has been becoming very challenging as the problem size and complexity increase tremendously, due to close couplings among network entities with heterogeneous service and resource requirements. By continuously interacting with the environment, deep reinforcement learning (DRL) provides a mechanism for different network entities to build knowledge and make autonomous decisions to improve network performance. In this article, we first review typical DRL approaches and recent enhancements. We then discuss the applications of DRL for mobile edge computing (MEC), which can be used for the low-power IoT devices, e.g., wireless sensors in healthcare monitoring, to offload computation workload to nearby MEC servers. To balance power consumption in offloading and computation, we propose a novel hybrid offloading model that exploits the complement operations of RF communications and low-power backscatter communications. The DRL framework is then customized to optimize the transmission scheduling and workload allocation in two communications technologies, which is shown to enhance the offloading performance significantly compared with existing schemes.

## Index Terms

Deep reinforcement learning, mobile edge computing, data offloading, backscatter communications.

Shimin Gong is with the School of Intelligent Systems Engineering, Sun Yat-sen University, China, email: gong-shm5@mail.sysu.edu.cn. Yutong Xie is with Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, email: yt.xie@siat.ac.cn. Jing Xu is with the School of Electronic Information and Communications, Huazhong University of Science and Technology, China, email: xujing@hust.edu.cn. Dusit Niyato is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore, email: dniyato@ntu.edu.sg. Ying-Chang Liang is with the Center for Intelligent Networking and Communications (CINC), University of Electronic Science and Technology of China, China, email: liangyc@ieee.org.

## I. INTRODUCTION

Modern wireless networks have to embrace the upsurge of traffic demands and diverse quality provisioning requirements. This requires a strategic shift in the network design that utilizes sophisticated wireless technologies in a more decentralized, ad-hoc, and diverse environment. As such, the network design problems become very challenging as the dimensionality and complexity rapidly increase, e.g., due to couplings among different network entities. Recently, deep reinforcement learning (DRL) has been developed as a breakthrough technology to learn the optimal control strategy in a dynamic network environment by continuously interacting with it [1]. DRL integrates deep neural networks (DNNs) with the conventional reinforcement learning algorithms for autonomous decision making. It becomes capable of solving high dimensional, non-convex, and model-free network control problems, e.g., channel access and resource allocation in mobile edge computing (MEC) [2]. These are very difficult to handle by conventional techniques such as convex optimization, dynamic and stochastic programming, due to imprecise modeling, uncertain system dynamics, and huge variable spaces. Hence, the application of DRL in wireless networks is envisioned to revolutionize the network optimization paradigm.

In this article, we first provide an overview of the DRL framework in Section II and its variants to improve the stability and learning performance. In Section III, as a concrete example, we shift our focus on performance optimization of the emerging MEC applications, which is generally complicated by the resource competition and interactions among multiple wireless users, base stations, caching and MEC servers [2]. We firstly build a general DRL framework to learn the optimal data offloading policy with uncertain network information, and then review the existing applications of DRL framework for MEC in different network scenarios. We observe that data offloading is not always preferred by low-power IoT devices due to the high energy consumption in wireless communications. Hence, in Section IV, we introduce a novel hybrid MEC offloading model to balance the energy consumption in offloading and computation. Besides local computation, the hybrid model allows data offloading to the MEC server via either the active RF communications or the passive wireless backscatter [3]. Our numerical results verify that the hybrid MEC offloading can significantly improve the network performance, by learning the optimal transmission scheduling and workload allocation among different offloading schemes. Finally, some open issues are discussed in Section V.

## II. AN OVERVIEW OF DEEP REINFORCEMENT LEARNING

In this section, we first review fundamentals of reinforcement learning and then discuss its extension to DRL, as well as various techniques to improve the learning efficiency and stability.

### A. Fundamentals of Reinforcement Learning

Reinforcement learning is an effective solution to Markov Decision Processes (MDPs), which is composed of the decision-making agent, system state, action, and reward [4]. The agent is the entity of decision making through interactions with the environment. Based on the observation of the environment, referred to as the system state, the agent takes an action and then receives an immediate reward correspondingly to the state-action pair. The action affects the environment and may cause the transition to a new system state. The immediate reward and the transition to a new state will guide, i.e., reinforce, the adaptation of the agent's policy, which defines the sequence of actions taken in each decision epoch as the system evolves. This learning process continues as we find the optimal policy to maximize the accumulated reward, which can be characterized by either the state-value or action-value function. The state-value records the expected total reward starting from an initial system state, while the action-value, also referred to as the  $Q$ -value, maps each state-action pair to the accumulated reward.

There are mainly value- and policy-based approaches for solving reinforcement learning problems [4]. The value-based approach estimates the value function and takes the action to improve it directly in an iterative process. The estimation of value functions can be based on value iteration following the Bellman equation or  $Q$ -learning algorithm. A variant of the value-based approach relies on the estimate of an advantage-value, which can stabilize the learning process by subtracting a baseline from the estimate of action-value. The policy-based approach improves the value function by updating a parametric policy in gradient-based methods. Reinforcement learning can be also classified into on- and off-policy approaches. The on-policy learning relies on the sample trajectory induced by the current policy, i.e., all future actions are chosen according to the current policy. This may require more interactions with the environment to ensure unbiased policy updates and thus make it impractical for solving complicated problems. The off-policy learning can improve the sample efficiency by utilizing all historical sample trajectories. However, it requires more effort in hyperparameter tuning to ensure the convergence in learning.

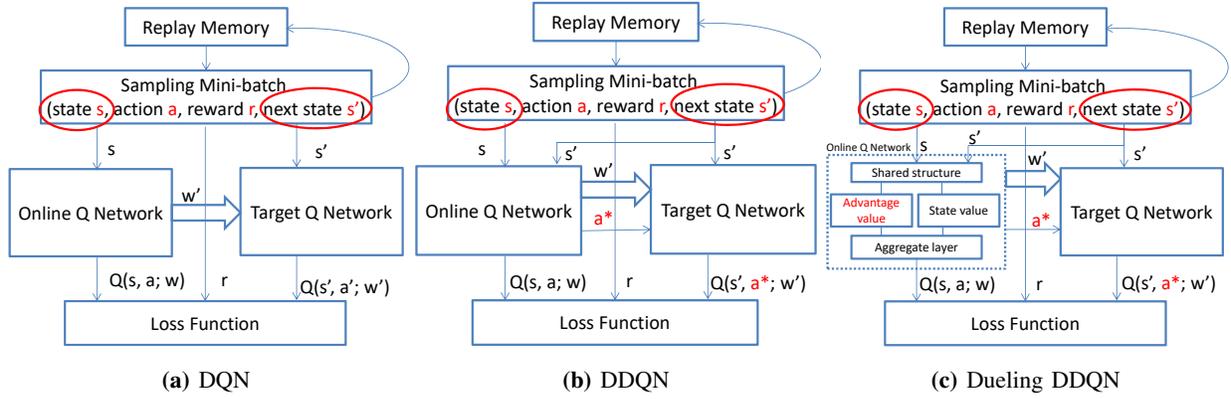
## B. Deep Reinforcement Learning Approaches

The reinforcement learning becomes unstable and even fail to converge when the state and action spaces are large in complex wireless networks. DRL can use the DNNs as function approximators for different components of reinforcement learning, including the value function, policy, and the underlying system model, e.g., the state transition probability. In the following, we introduce the basics of DRL and recent advances to improve its learning performance.

1) *Deep Q-Network (DQN)*: It extends the value-based  $Q$ -learning algorithm for MDPs by using DNN as a parametric approximation for the action-value function [5]. The success of DQN and its variants relies on two key mechanisms, i.e., experience replay and target  $Q$ -network, to stabilize learning with large state and action spaces. The experience replay maintains a replay memory to buffer historical transition samples and randomly selects a subset of samples, i.e., mini-batch, to train the DNN. This can break the sample correlations and ensure more efficient training by independent samples. The training of DNN in each step aims to minimize the temporal-difference (TD) error, i.e., the mean-squared difference between the estimated  $Q$ -value by the DNN and its target value. Practically, we can replay more frequently the transition samples that generate a higher expected reward. Hence, a prioritized experience replay (PER) scheme can potentially increase the learning speed [6]. A straightforward way for PER is to prioritize samples by their TD-errors. A higher TD-error implies a larger potential to be further optimized. The TD-error based PER can be further combined with random sampling to ensure that all transition samples have the chance to be selected for training [4].

2) *Double and Dueling DQN*: DQN uses a separate  $Q$ -network to generate the target value. The target  $Q$ -network updates its parameter by copying it from the online  $Q$ -network in every a few steps, as illustrated in Fig. 3(a). This can make the learning more stable compared to the  $Q$ -learning algorithm. The drawback of DQN lies in that it uses the  $\epsilon$ -greedy policy to select an action and evaluate it by the same  $Q$ -network [4]. This may lead to over-optimistic estimation of the  $Q$ -value. To correct this, Double DQN (DDQN) updates the action by the online  $Q$ -network and then evaluates it by the target  $Q$ -network [7], as illustrated in Fig. 3(b). Another variant of DQN decomposes the  $Q$ -value into two streams, i.e., the state-value and the advantage-value [8], approximated by two independent DNNs in a dueling architecture. The two streams are then combined via an aggregating layer to produce the final estimate of the  $Q$ -value.

3) *Deep Deterministic Policy Gradient*: The policy-based and value-based approaches can be combined in the actor-critic framework [4]. The critic function produces the estimation of the



DRL approaches	Year	Action/State	DNN setting	Learning speed	Sample efficiency	On/Off-policy
Plain DQN [5]	2015	Disc. moderate	One DNN	Slow	Low	Off-policy
DQN with PER [6]	2016	Disc. moderate	One DNN	Moderate	High	Off-policy
DDQN with PER [7]	2016	Disc. large	Two DNNs	Fast	High	Off-policy
Dueling DDQN [8]	2016	Disc. large	Three DNNs	Fast	High	Off-policy
DDPG [9]	2016	Cont. large	Two DNNs	Fast	Moderate	Off-policy
Rainbow [10]	2018	Disc. large	Adaptive	Faster	High	Mixed
TRPO [11]	2019	Cont. large	Two DNNs	Fast	Moderate	On-policy

(d) Comparison of different DRL approaches.

**Fig. 1:** The comparison of typical value-based and policy-based DRL approaches.

$Q$ -value by minimizing the TD-error. The actor function then updates the policy parameter using the critic's feedback. Two independent DNNs can be used as the parametric approximations for the critic and actor functions, respectively. The intuition behind actor-critic framework stems from the policy gradient theorem that builds the connection between policy gradient and the  $Q$ -value. It decomposes the gradient computation into the evaluation of the  $Q$ -value and the gradient of the parametric policy, averaged over the whole state and action spaces. One recent development is to extend the policy gradient theorem to deterministic policy gradient (DPG), which outputs a deterministic action instead of a distribution on action space and thus makes it more efficient to estimate the policy gradient. The deep deterministic policy gradient (DDPG) algorithm combines DQN and DPG in the actor-critic framework to make the learning more stable and robust by using the experience replay and target  $Q$ -network for DNN training [9].

A comparison of typical DRL approaches is listed in Fig. 3. In general, DQN and its variants are applicable to discrete action space, which are natural extensions of  $Q$ -learning algorithm

for solving MDPs with large action and state spaces. The Rainbow algorithm in [10] is an integrated design of different DQN variants, which achieves the best learning speed and maximum reward. The continuous action space can be more preferably tackled by DDPG in [9] and the trust-region policy optimization (TRPO) in [11]. To avoid large variance in gradient estimation, TRPO formulates a constrained optimization to search for a better policy that improves the value function. Besides, we observe that the off-policy is more popular for DRL as it can use all historical samples efficiently. Though TRPO is generally on-policy, it has been adapted in [11] to leverage a replay buffer and thus can achieve a better learning performance compared to DDPG.

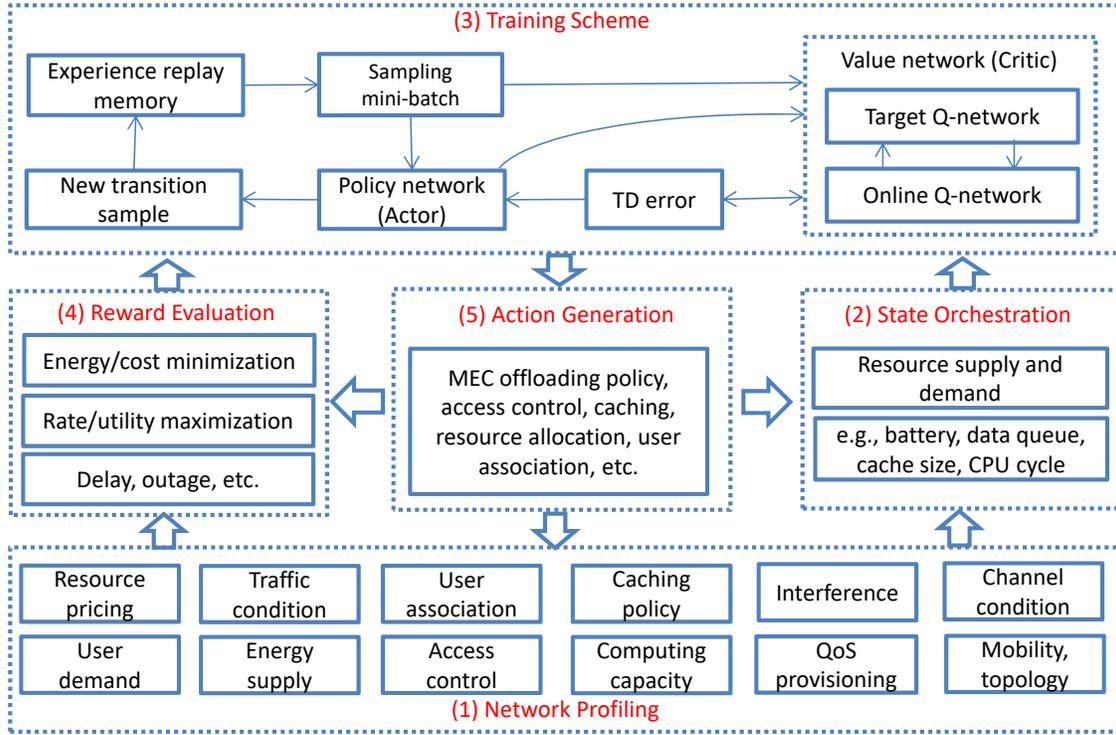
As modern wireless networks become large-scale and complicated, the network control problems face very diversified decision variables, including both discrete indicators and continuous variables for resource allocation. Thus, both value- and policy-based methods need to be used jointly for mixed decision-making problems. In the following, we focus on the applications of DRL in the emerging MEC offloading scenarios, which typically involve the interactions among multiple network entities and complicated optimization in both discrete and continuous domains.

### III. DRL-BASED DATA OFFLOADING FOR MOBILE EDGE COMPUTING

MEC offloading allows IoT devices to offload data and computation-intensive workload (e.g., compressing and encryption) to resource-rich MEC servers. It can potentially reduce the processing delay, extend the battery lifetime, and even enhance security for IoT applications [2]. One of the critical design issues is to optimize the offloading rate, workload allocation, and choose the optimal MEC server, considering the time-varying channel conditions, user mobility, energy supply, dynamic workloads, and various resource constraints. A joint optimization on caching, offloading, networking, and transmission control is usually very complicated due to close couplings among multiple wireless users, base stations, and MEC servers. The optimization is also very inflexible to capture the network dynamics with uncertain parameters, e.g., the fluctuating channel conditions, the time-varying workload and energy supplies.

#### A. General DRL Framework for MEC Offloading

DRL avoids above-mentioned difficulties by reformulating the network control problem into the MDP framework and enhancing the reinforcement learning solution by deep learning techniques. In the sequel, we propose a general DRL framework for MEC offloading that can be



**Fig. 2:** DRL-based MEC offloading framework.

flexibly tailored for learning offloading strategy under different network scenarios. As illustrated in Fig. 2, the DRL framework includes the following main components:

- (1) **Network Profiling:** The network environment contains very high dimensional information. Dimension reduction is required to speed up the learning process. Network profiling helps to extract problem-dependent information closely related to the network control problems. This can be assisted by conventional model-based optimization problems.
- (2) **State Orchestration:** It aims to select the most salient or indicative state variables to minimize the state space without compromising the learning performance. The network performance depends on the demands and supplies of various resources. Hence, the system state can be set to show the real-time dynamics of resource consumption and its regeneration.
- (3) **Training Scheme:** The training scheme can flexibly organize the value and policy networks to learn both discrete and continuous offloading decisions, e.g., the discrete indicators for base station (de)activation, channel assignment, user association, and routing selection, as well as the continuous variables for bandwidth allocation and beamforming optimization.
- (4) **Reward Evaluation:** The reward in each decision epoch drives the DRL agent to adjust its MEC offloading policy. Practically, the reward is evaluated after completing the workload

after a few decision epoches or time slots. A model-based optimization can be deployed to estimate the instant reward based on the prediction of future network dynamics.

- (5) *Action Generation*: The DRL agent outputs a vector of actions for each system state, which will be translated into the control variables to execute the offloading decisions. Quantization or approximation can be required in some cases to project continuous variables into discrete actions. Random noise can be also added to continuous actions for a better exploration.

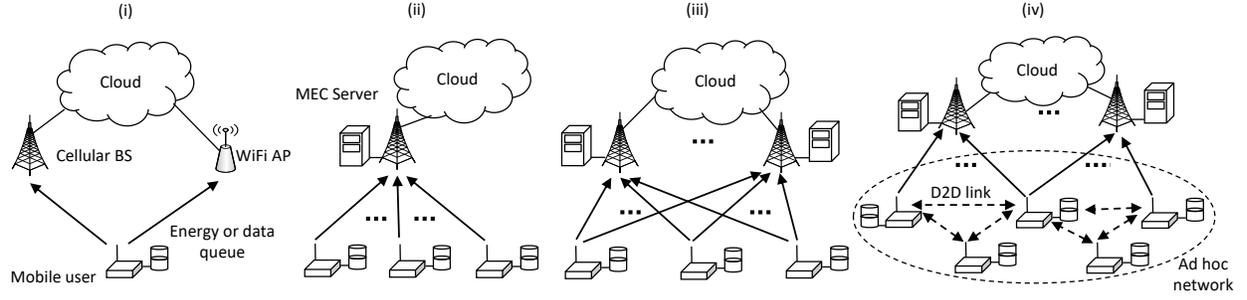
The general DRL framework can be applied to optimize MEC offloading policies under different network scenarios by customizing different components of the DRL framework to meet the performance requirements of various design problems. In the following, we provide a review of recent works on the applications of DRL for MEC offloading problems.

### *B. Design Issues for DRL-based MEC Offloading*

1) *Network Selection for Cost Minimization*: In the simplest case with one wireless user and multiple access points, e.g., cellular base station and WLAN access point [12], the MEC offloading is regarded as a network selection problem as illustrated in case (i) of Fig. 3. The wireless user can either access the cellular network or WLAN with different costs. To minimize the user's energy consumption and cost for channel access, DQN can be constructed to learn the optimal selection scheme without knowing the user's mobility pattern. The offloading decision is made based on the prediction of the user's location and the remaining data size.

2) *Channel and Capacity Sharing*: When multiple wireless users request for the computation resources simultaneously from a single MEC server, e.g., [13], as shown in case (ii) of Fig. 3, the spectrum and capacity sharing becomes a critical problem to minimize the cost of delay and power consumptions for all users. The system state can be the sum cost of all users and the remaining capacity of the MEC server. The DRL agent learns the continuous resource allocation for wireless users and the binary offloading decisions, considering a limited capacity of the MEC server and time-varying channel conditions.

3) *MEC Server and User Association*: With multiple base stations or MEC servers, each wireless user's computation offloading can be routed via different base stations, as shown in case (iii) of Fig. 3. To minimize the cost of processing delay, the authors in [14] employed DDQN to learn the optimal offloading policy including the binary user association and transmit control strategies. The system state consists of the channel conditions between the wireless users and the base stations, the statuses of energy and data queues. Considering low utilization of base



(a) Applications of DRL-based MEC offloading framework in different network scenarios: (i) Network selection [12], (ii) Channel and capacity sharing [13], (iii) MEC server and user association [14], (iv) Collaborative data offloading [15].

Scenarios	DRL models	States	Actions	Rewards
(i) [12]	DQN using CNN	User's location and the remaining file size	Keep idle, access WLAN or cellular network	Minimize cost and energy consumption
(ii) [13]	DDPG	Channel states and all users' task buffers	Power distribution in local and offloading computation	Minimize energy consumption within delay deadline
(iii) [14]	Double DQN	Channel states, energy and task queues	User association and power allocation	Maximize satisfaction on delay, outage probability, and payment for MEC service
(iv) [15]	Plain DQN	Task queues of all users and their interdistance	Task distribution in local and offloading computing	Maximize utility minus the cost of energy consumption, delay, and outage probability

(b) The states, actions, rewards definitions in different DRL approaches for MEC offloading.

**Fig. 3:** DRL applications in different MEC offloading scenarios.

stations, DDQN can be also used to control the (de)activations of base stations to reduce total energy consumption while maintaining the same quality provisioning.

4) *Collaborative Data Offloading*: Besides offloading to the MEC server, the collaborative offloading among multiple wireless users can be envisioned in case (iv) of Fig. 3, i.e., each wireless user can offload its computation workload to nearby users via device-to-device communications, e.g., [15]. The optimization of offloading decisions depends on the number of remaining tasks at each wireless user, the availability of the computation resources, and the link quality between wireless users. DQN or DDQN can be customized to learn the optimal offloading policy in a mobile ad-hoc network to maximize the resource utilization or minimize total power consumption, subject to the user's energy and delay requirements.

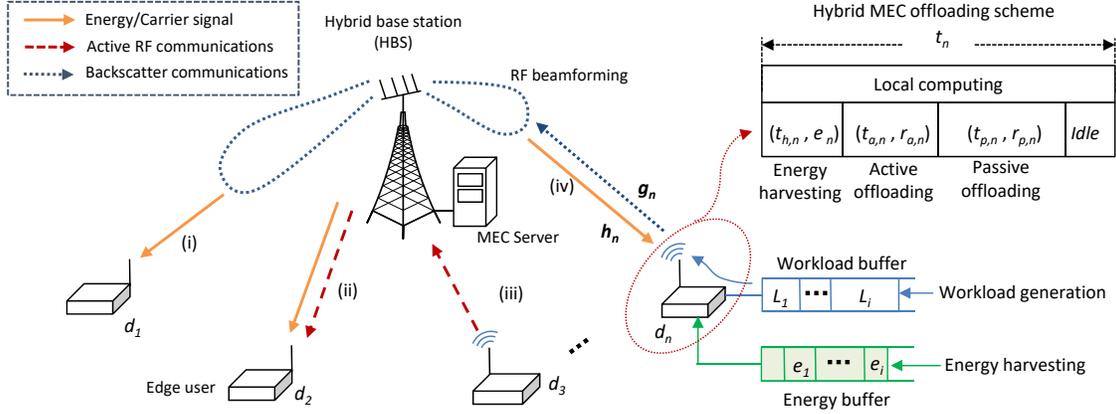
#### IV. DRL APPROACH FOR A HYBRID MEC OFFLOADING MODEL

One design objective of the future wireless network is to embrace the ubiquitous interconnections of low-power IoT devices, e.g., the wearable wireless sensors for healthcare monitoring, either battery-powered or wireless powered via energy harvesting [13]. For these low-power IoT devices, it is clear that energy consumption on data processing can be reduced significantly by offloading computation-intensive workload to the MEC servers, e.g., [12]–[15]. However, in another aspect, the energy saving on computation comes with the price of more energy consumption on computation offloading, which is generally performed by conventional RF communications. Due to the high energy consumption of RF communications, MEC offloading may not be affordable by these low-power IoT devices.

In this section, we tackle this problem by proposing a hybrid offloading strategy that can schedule data offloading in both high-rate RF communications and low-power backscatter communications [3]. The backscatter radio operates in the passive mode by reflecting the incident RF signals. It is featured with extremely low power consumption and a low data rate, while the active radio in RF communications can transmit reliably with a higher data rate by adapting its transmissions against the channel fading effect. Hence, we aim to optimize the hybrid MEC offloading policy to balance energy consumptions in both data offloading and computation. This can be achieved by exploiting the complement operations of the passive and active radios. However, due to the couplings among two radio modes, it becomes more complicated to optimize the MEC offloading policy by using the conventional model-based approaches. In the sequel, we employ the DRL framework to optimize the hybrid MEC offloading strategy with uncertain channel conditions, dynamic energy supply, and time-varying workloads at the IoT devices.

##### A. Hybrid MEC Offloading Model

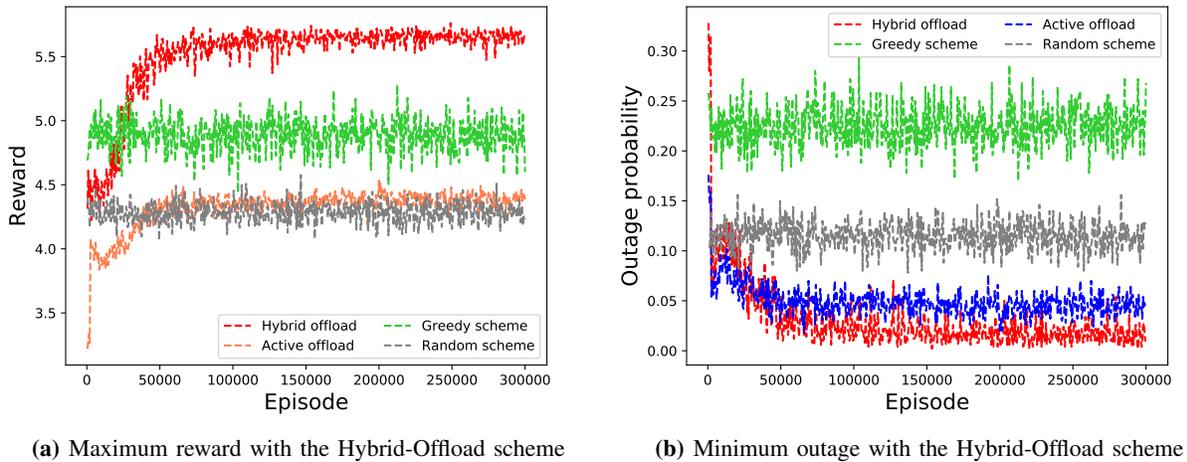
We consider a set of edge users, e.g., wireless IoT sensors, that send backlogged workloads to a hybrid base station (HBS), which is co-located with the MEC server. The system model is illustrated in Fig. 4. The channel from the HBS to each edge user is modeled by a finite-state Markov model. The HBS allocates each edge user a time slot for MEC offloading, similar to the time-slotted structure in [13]. Each edge user can harvest RF energy from the HBS and the ambient RF sources with random power density. The energy harvesting capability is illustrated in case (i) of Fig. 4. The edge user's workload is uncertain due to the user's mobility and time-varying demand of upper layer applications, e.g., the data sampling rate may vary according to



**Fig. 4:** Hybrid MEC offloading scheme: (i) RF-based energy harvesting, (ii) Active offloading via RF communications, (iii) Down-link power and information transfer, (iv) Passive offloading via backscatterer communications.

the health conditions of the subject being monitored. The user's workload needs to be processed locally or remotely at the MEC server before a time deadline. We assume that the MEC server can return the processed data to the edge user instantly via simultaneous power and information transfer, as illustrated in case (ii) of Fig. 4. The hybrid MEC offloading scheme allows each user to flexibly switch data offloading between the passive backscatter communications and the active RF communications, as illustrated in cases (iii) and (iv) of Fig. 4, respectively. To maintain a fixed offloading rate, the active radio's transmit power has to be adapted with the time-varying channel conditions. This implies a dynamic process of the edge user's energy buffer.

It is obvious that the offloading scheduling between two radio modes introduces an additional degree of freedom to improve the MEC performance in such a dynamic network environment. The DRL approach aims to learn the optimal hybrid MEC offloading policy from past experience. Given the channel conditions, energy status, and workload in each time slot, the edge user will choose its action (e.g., local computation, passive or active offloading) to maximize the reward function, which is defined as the energy efficiency, i.e., the successfully processed workload per unit energy. Workload outage happens when the edge user's workload is not processed successfully before the delay bound. In this case, the instant reward will be set to zero. To proceed, we divide each time slot into flexible sub-slots as illustrated in Fig. 4. The first sub-slot  $t_{h,n}$  is reserved for RF energy harvesting. The following sub-slot  $t_{a,n}$  is allocated to active offloading with a higher rate  $r_{a,n}$  and another sub-slot  $t_{p,n}$  is used by passive offloading with a lower rate  $r_{p,n}$ . The offloading schemes also differ in their power consumption. To achieve the maximum energy efficiency, the DRL agent is designed to learn a transmission scheduling



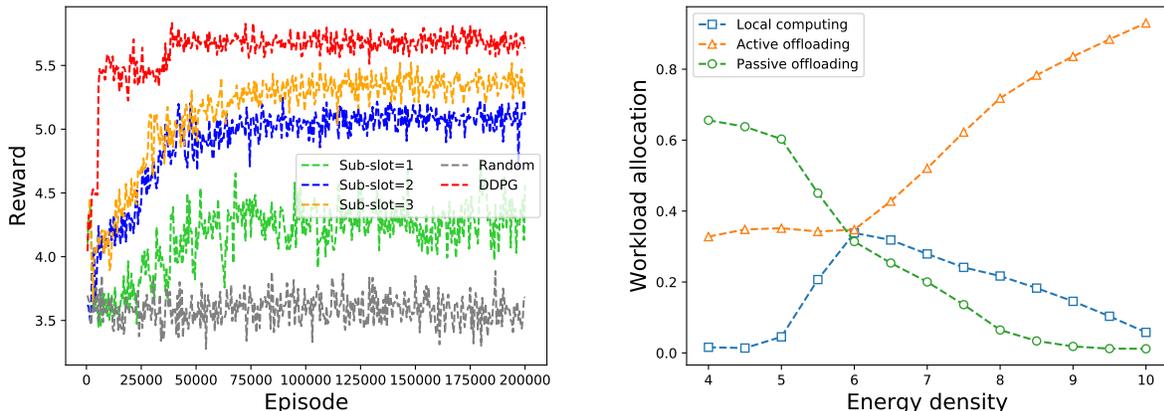
**Fig. 5:** Performance comparison among different MEC offloading schemes.

policy that determines the optimal action on each system state, including time and workload allocations among local computation, active and passive offloading, subject to the edge user’s energy budget constraint.

### B. Numerical Evaluation

To exploit the flexibility and performance gain via hybrid MEC offloading scheme, we compare it with the conventional offloading scheme, namely, the Active-Offload scheme that only supports active RF communications, e.g., [13]. We also compare it with the typical greedy and random schemes. In Fig. 5, we show the performance of different schemes and observe that the DRL-based Hybrid-Offload scheme achieves the highest reward and the lowest outage performance, as shown in Fig. 5(a) and Fig. 5(b), respectively. The Active-Offload scheme uses a similar DRL framework as that of the Hybrid-Offload scheme, however with the reduced action space, i.e., it only chooses between local computation and active offloading. Hence, it achieves a reduced reward performance than that of the Hybrid-Offload scheme. The benchmark greedy scheme always chooses the myopic action to maximize the instant reward in each time slot. It even performs better than the Active-Offload scheme due to its flexibility in radio mode switching.

In the next group of simulations, we equally divide each time slot into multiple sub-slots and assume that the edge user follows the same DRL framework to optimize its offloading decision independently in each sub-slot. By this way, we can flexibly allocate the workload and thus approximate the optimal workload allocation strategy among local computation, passive



(a) Higher reward can be achieved with more sub-slots

(b) Workload allocation among three computation schemes

**Fig. 6:** Performance comparison with different parameter settings.

and active offloading. In Fig. 6(a), we show the performance of the Hybrid-Offload scheme when we set a different number of sub-slots for MEC offloading. We observe that it generally achieves a higher reward performance with more sub-slots. We also show the performance of the DDPG algorithm for continuous control that is shown to achieve the maximum reward. Such a performance gain is obtained from the increased flexibility in workload and time allocation.

In Fig. 6(b), we show the averaged workload allocation among different computation schemes at the convergence of the Hybrid-Offload algorithm. The  $x$ -axis of Fig. 6(b) denotes the mean power density in the ambient RF environment. We can observe that with low energy supply the passive offloading scheme is preferred due to its extremely low power consumption. With a higher energy density, the edge user generally harvests more RF energy and thus it prefers to use the active offloading scheme. This can provide a higher offloading rate and thus reduce the processing delay.

## V. OPEN RESEARCH ISSUES

Though DRL has been successfully applied to various network control problems, there still exist some challenges and open issues for MEC offloading in wireless networks.

1) *Multi-agent DRL for MEC Offloading:* MEC offloading involves multiple heterogeneous network entities, e.g., wireless users, base station, and MEC server, which may have totally different reward functions and control variables. Each user can customize its own DRL framework

and make decisions based on local observations. However, this may destroy the Markovian property of the underlying system model and lead to divergent learning performance.

2) *Model-based Reward Evaluation*: The DRL agent requires real-time reward evaluation to drive the learning process. As the performance of MEC offloading decision is usually not observable until the completion of workload, we require a more effective way combining learning and model-based optimization to predict the reward with incomplete network information.

3) *Hierarchical DRL for MEC Offloading*: MEC offloading decision generally involves both discrete and continuous control variables. To improve learning efficiency, a hierarchical or two-stage DRL framework can be implemented to learn the optimal resource allocation strategy by using the policy-based DRL approaches in the inner loop, and then update the discrete user association or offloading decisions by DQN or its variants in the outer loop.

## VI. CONCLUSIONS

In this paper, we firstly have reviewed the DRL framework for its applications in MEC offloading with uncertain network information. Then, we have customized the DRL framework to realize a novel hybrid MEC offloading scheme that exploits the complement transmissions of the passive and active radios. Numerical results demonstrate that it can significantly improve the offloading performance. In the last, we have outlined a few open research issues.

## REFERENCES

- [1] N. C. Luong, D. T. Hoang, S. Gong, D. Niyato, P. Wang, Y. Liang, and D. I. Kim, "Applications of deep reinforcement learning in communications and networking: A survey," *IEEE Commun. Surv. Tut.*, May 2019.
- [2] N. Abbas, Y. Zhang, A. Taherkordi, and T. Skeie, "Mobile edge computing: A survey," *IEEE Internet of Things Journal*, vol. 5, no. 1, pp. 450–465, Feb. 2018.
- [3] S. Gong, J. Xu, D. Niyato, X. Huang, and Z. Han, "Backscatter-aided cooperative relay communications in wireless-powered hybrid radio networks," *IEEE Network*, May 2019.
- [4] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. The MIT Press, 2018.
- [5] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
- [6] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," in *Proc. Int. Conf. Learning Representations*, Puerto Rico, May 2016.
- [7] H. v. Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double Q-learning," in *Proc. AAAI Conf. Artificial Intelligence*, Phoenix, Arizona, USA, Feb. 2016, pp. 2094–2100.
- [8] Z. Wang, T. Schaul, M. Hessel, H. Van Hasselt, M. Lanctot, and N. De Freitas, "Dueling network architectures for deep reinforcement learning," in *Proc. Int. Conf. Machine Learning (ICML)*, New York, NY, USA, Jun. 2016, pp. 1995–2003.

- [9] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” in *Proc. Int. Conf. Learning Representations (ICLR)*, San Juan, Puerto Rico, May 2016.
- [10] M. Hessel, J. Modayil, H. P. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. G. Azar, and D. Silver, “Rainbow: Combining improvements in deep reinforcement learning,” in *Proc. AAAI Conf. Artificial Intelligence*, 2018.
- [11] D. Kangin and N. Pugeault, “On-policy trust region policy optimisation with replay buffers,” in *Proc. Int. Conf. Learning Representations (ICLR)*, New Orleans, Louisiana, USA, May 2019.
- [12] C. Zhang, Z. Liu, B. Gu, K. Yamori, and Y. Tanaka, “A deep reinforcement learning based approach for cost-and energy-aware multi-flow mobile data offloading,” *IEICE Trans. Commun.*, vol. E101-B, no. 7, pp. 2017–2025, 2018.
- [13] Z. Chen and X. Wang, “Decentralized computation offloading for multi-user mobile edge computing: A deep reinforcement learning approach,” *CoRR*, vol. abs/1812.07394, Dec. 2018.
- [14] X. Chen, H. Zhang, C. Wu, S. Mao, Y. Ji, and M. Bennis, “Optimized computation offloading performance in virtual edge computing systems via deep reinforcement learning,” *IEEE IoT Journal*, vol. 6, no. 3, pp. 4005–4018, Jun. 2019.
- [15] D. V. Le and C.-K. Tham, “A deep reinforcement learning based offloading scheme in ad-hoc mobile clouds,” in *Proc. IEEE INFOCOM (Workshop)*, Honolulu, USA., Apr. 2018.