# Privacy-Preserved Task Offloading in Mobile Blockchain with Deep Reinforcement Learning

Dinh C. Nguyen, *Student Member, IEEE,* Pubudu N. Pathirana, *Senior Member, IEEE,*
Ming Ding, *Senior Member, IEEE,* Aruna Seneviratne, *Senior Member, IEEE*

*Abstract*—Blockchain technology with its secure, transparent and decentralized nature has been recently employed in many mobile applications. However, the process of executing extensive tasks such as computation-intensive data applications and blockchain mining requires high computational and storage capability of mobile devices, which would hinder blockchain applications in mobile systems. To meet this challenge, we propose a mobile edge computing (MEC) based blockchain network where multi-mobile users (MUs) act as miners to offload their data processing tasks and mining tasks to a nearby MEC server via wireless channels. Specially, we formulate task offloading, user privacy preservation and mining profit as a joint optimization problem which is modelled as a Markov decision process, where our objective is to minimize the long-term system offloading utility and maximize the privacy levels for all blockchain users. We first propose a reinforcement learning (RL)-based offloading scheme which enables MUs to make optimal offloading decisions based on blockchain transaction states, wireless channel qualities between MUs and MEC server and user's power hash states. To further improve the offloading performances for larger-scale blockchain scenarios, we then develop a deep RL algorithm by using deep Q-network which can efficiently solve large state space without any prior knowledge of the system dynamics. Experiment and simulation results show that the proposed RL-based offloading schemes significantly enhance user privacy, and reduce the energy consumption as well as computation latency with minimum offloading costs in comparison with the benchmark offloading schemes.

*Index Terms*—Blockchain, mobile edge computing, mining, task offloading, privacy, deep reinforcement learning.

## I. INTRODUCTION

In recent years, the blockchain technology has been employed widely in various industrial applications such as Internet of Things (IoT), healthcare, industrial applications, etc. [1], [2]. Blockchain works as a peer-to-peer public ledger where users can store data (i.e., records of transactions) and share information with other blockchain nodes in a trustworthy and decentralized manner. With the advancement of mobile technologies, blockchain now can be implemented in mobile devices to provide more flexible blockchain-based solutions

for IoT applications [3], [4]. The foundation of the efficient and secure operation of blockchain is a computation process known as *mining* [5]. In order to append a new transaction to the blockchain, a blockchain user, or a miner, needs to run a mining puzzle, i.e., Proof of Work (PoW) or Proof of Stake (PoS) which is generally complicated and requires intensive computations. Resource-limited IoT nodes or mobile devices, therefore, cannot participate in the mining operation, which can restrict the application of blockchain in mobile systems. Recently, mobile edge computing (MEC), which enables mobile devices to offload their computation tasks to a nearby computationally powerful MEC server, provides effective solutions to solve computation issues of mobile devices [6]–[8]. By employing highly computational MEC resources, the performance of MEC-based mobile systems may be improved significantly, such as reducing computation latency and improving computation efficiency [9]–[11].

In the literature, to solve the offloading issues in mobile systems, various approaches have been proposed, using convex optimization-based methods [12], [13], or machine learning techniques such as reinforcement learning (RL) [14]–[16] or deep reinforcement learning (DRL) [17]–[20]. In mobile blockchain applications, MEC-based computation offloading strategies to optimize mining services [21] have been considered in some previous studies [22]–[24]. Such pioneering studies consider offloading issues, from data offloading, resource allocation and block size optimization to enhance offloading efficiency in blockchain networks. Addition to computation offloading, privacy is also another critical issue that should be considered carefully when designing offloading systems. Offloading data tasks and mining tasks to MEC server to reduce computing burden on mobile devices while ensuring high privacy for mobile users is very important in network management and ensures the robustness of the blockchain-based systems. However, we find that the study that can achieve optimal policies of offloading mining tasks with enhancing privacy awareness in mobile blockchain is still missing. Many pioneering blockchain schemes [22]–[24] have concentrated on only offloading of mining tasks without considering user privacy preservation, which has been one of the key challenges in MEC-based blockchain networks.

To this end, we propose a DRL-based dynamic task offloading scheme for a MEC blockchain network where each mobile user (MU) acts as a miner that can offload its IoT data processing tasks and mining tasks to the MEC server for computing services. Particularly, we formulate task offloading and user privacy preservation as a joint optimization problem.

Then we develop a RL-based algorithm to solve the proposed optimization problem with a simplified blockchain model. Next, to break the curse of high dimensionality in state space when increasing the number of blockchain users, a DRL-based method called deep Q-network (DQN) algorithm is proposed. The objective of the proposed scheme is to obtain optimal offloading actions for all blockchain miners such that the user privacy level is maximized while minimizing computation latency and energy consumption for offloading. It is worth noting that in this paper, we only concern about the optimization of the offloading decision of each user, while edge resource is assumed to be sufficient to perform all offloading tasks. The main differences between our study and the existing schemes [21]–[24] are highlighted as follows.

1) We consider a new task offloading model for both IoT data processing tasks and mining tasks in a MEC-enabled blockchain network. Mobile users can act as miners to learn dynamic computation offloading policies to perform IoT data processing and participate in the mining puzzle to make extra profits.

2) To obtain optimal offloading policies for all miners, we propose a RL-based algorithm and then use a deep Q network to implement our proposed offloading strategies, aiming to achieve the best privacy performance and mining profits while minimizing offloading latency and energy cost.

3) We conduct both experiments and numerical simulations to verify the advantages of the proposed offloading algorithm over the other offloading schemes.

The remainder of this paper is organized as follows. Section II reviews the related MEC-blockchain works. Section III introduces the system model, and then we present the formulation of computation model in Section IV. Based on system formulation, in Section V we derive the offloading optimization problem with a RL-based solution. Next, we propose two offloading schemes in Section VI by using RL and DRL. Experiment and simulation results are given in Section VII. Finally, our conclusions are drawn in Section VIII.

## II. RELATED WORKS

The works in [12], [13] were introduced with the objective of minimizing the energy usage or task execution latency for edge task offloading by leveraging traditional convex optimization methods, but they usually require prior system knowledge. RL has emerged as a strong alternative that can derive an optimal solution via trial and error without requiring any prior environment information [14]. Nevertheless, in high-dimensional offloading problems, the dimension of state and action space can be extremely high that makes RL-based solutions inefficient [15], [16]. Fortunately, DRL methods [17] such as DQN have been introduced as a strong alternative to deal with such high-dimensional problems and demonstrates its scalability and offloading efficiency in various MEC-based applications, such as multi-base station virtual MEC [18], multi-user MEC system [19], and multi-IoT networks [20], [21].

Further, a computation offloading scheme was proposed in [22] for a wireless blockchain networks where mobile users can offload their blockchain mining tasks to a MEC server or to a network of device-to-device users, subject to probabilistic backhaul and latency constraints. Meanwhile, the problem of computation power allocation for mining task offloading in the MEC-enabled blockchain was investigated in [23] with the main focus on reward optimization for mobile terminals. Further, the work [24] considered an optimization problem of offloading scheduling, resource allocation and adaptive block size scheme for a blockchain-based video streaming system with MEC. The user targets to offload computation to MEC nodes or mobile users such that transcoding revenue is maximized while achieving the high latency efficiency.

In addition to mining efficiency, privacy of blockchain users is another important issue that needs to be considered for mobile offloading in MEC-based blockchain networks [25]–[27]. However, most existing computation offloading frameworks for blockchain mining services [22]–[24] have ignored user privacy. A user privacy model was proposed for MEC-based mobile networks where mobile devices can select efficient offloading decisions using a constrained Markov decision process (CMDP) [28]. Meanwhile, the work in [29] proposed a privacy-aware computation offloading scheme using reinforcement learning, which enables IoT devices to learn an offloading policy so as to protect their personal information while achieving minimum system utility loss.

## III. SYSTEM MODEL

### A. Blockchain Network

In this paper, we consider a public blockchain network for MEC-enabled IoTs as shown in Fig. 1. The blockchain network consists of the MEC server, a wireless access point (AP), a set of multiple IoT devices and $N$ MUs denoted by a set $\mathcal{N} = \{1, 2, ..., N\}$. Blockchain is used as decentralized database to record IoT transactions and manage securely data exchange between IoT devices, MUs and MEC server. In our blockchain scenario, each IoT device is regarded as a lightweight node which is responsible to perform data collection from external sensing environments (e.g., wearable sensors for sensing health data [2]) and transmit wirelessly such data to MUs for necessary data processing. It is assumed that IoT devices are grouped and managed by a certain MU in a local private network. Specially, MUs represent their IoT devices to perform transactions with the MEC-blockchain network and execute sensing data. Here, each MU performs two computation functionalities for: IoT data processing tasks and mining tasks. For IoT data processing tasks, MUs can decide to execute locally or offload tasks to the MEC server for efficient computation. Meanwhile, for mining tasks, MUs can act as miners and purchase hash power from the edge cloud providers to perform mining and make extra profits [30]. (MU, miner and mobile miner terms are used interchangeably throughout the paper). In this way, miners can gather transactions of the blockchain network and group them into data blocks, which are then verified and appended into the blockchain after consensus verification. The miner which is the first one to solve the consensus puzzle is rewarded for their mining contribution, and the verified transaction is broadcast to other blockchain nodes in secure and immutable ways.
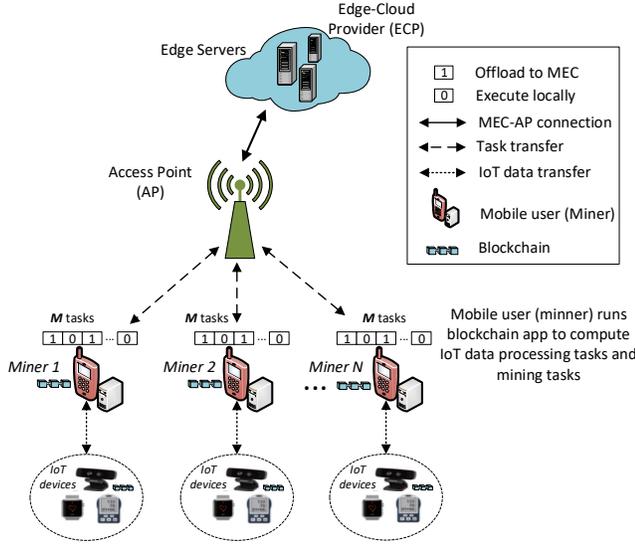
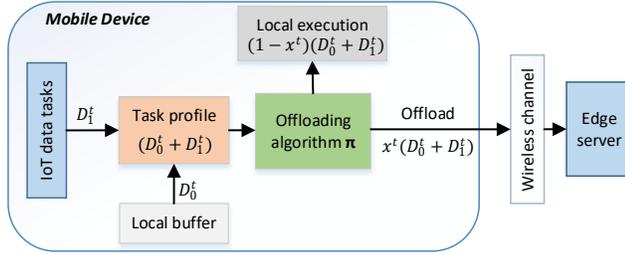Fig. 1: The proposed mobile edge blockchain architecture.



Fig. 2: The offloading procedure in mobile edge blockchain.

TABLE I: List of notations.

| Notation | Description |
|----------|-------------|
| $D_{1nm}^t$ | Blockchain transaction for task $m$ at miner $n$ at time slot $t$ |
| $D_{0nm}^t$ | Blockchain transaction in the buffer of miner |
| $X_{nm}^t$ | The total number of CPU cycles per task |
| $\tau_{nm}^t$ | Completion deadline for mining a task |
| $x_{nm}^t$ | The task offloading policy |
| $g^t$ | Wireless channel power gain |
| $f_{nm}$ | MEC computational capacity |
| $L_{nm}^{o,t}$ | Latency of edge computation |
| $E_{nm}^{o,t}$ | Energy consumption of edge computation |
| $L_{nm}^{l,t}$ | Latency of local computation |
| $E_{nm}^{l,t}$ | Energy consumption of local computation |
| $P_{nm}^t$ | The total task offloading privacy level |
| $p_n$ | The miner $n$'s hash power purchased from edge-cloud provider |
| $R_n^{mining}$ | The mining reward of miner $n$ |
| $\alpha$ | The learning rate of deep Q-learning |
| $\gamma$ | The discount factor of deep Q-learning |
| $\mathcal{D}$ | The replay memory capacity |

the current IoT data into two parts: data processing task for execution and the rest for storage in the buffer for future process. The mining strategy with task offloading is shown in Fig. 2 and can be explained as follows.

At each time slot $t$, the miner receives new IoT data $D_{1nm}^t$ from IoT devices and adds it into a block based on the blockchain concept [1]. Besides, the miner also has to process the previous IoT data $D_{0nm}^t$ stored in the buffer. Moreover, for simplicity, we assume that each block contains only a transaction which will be verified by the miner through a mining process (e.g., PoW) so that the verified transaction can be appended to blockchain. As a result, the total blockchain data of a transaction at each time slot $t$ is $(D_{0nm}^t + D_{1nm}^t)$ which can be formulated as a data processing task $R_{nm}^t$. The miner will choose to offload the data processing task to the MEC server or execute it locally and store the rest data processing tasks into the local buffer for processing in the future. Accordingly, we consider two computing modes in this paper, namely:

- *Offloading computing*: The miner $n$ offloads its data processing tasks to the nearby MEC server.
- *Local computing*: The miner $n$ executes its data processing tasks at the local mobile devices.

Specially, we introduce an offloading decision policy denoted by a binary variable $x_{nm}^t \in \{0,1\}$. Here $x_{nm}^t = 1$ means that the miner $n$ offloads the task $m$ to the MEC server over the wireless channel with the wireless channel power gain $g(t)$, and $x_{nm}^t = 0$ means that the miner $n$ decides to process the task $m$ locally. For simplicity, the wireless channel power gain $\{g^t\}_{t \geq 0}$ can be formulated as a Markov chain model with $Pr(g^{t+1} = m | g^t = n) = h_{nm}, \forall n, m \in \mathcal{G}$, where $\mathcal{G}$ is the wireless channel state set. Particularly, we define that $g^t = 1$ or $(g^t = 0)$ represents the good (or bad) channel state at the timeslot $t$. This policy will be used to analyze channel conditions for ensuring user privacy during task offloading. The notations used in this paper are summarized in Table I.

## C. Offloading For Blockchain Mining Tasks

In addition to IoT data computation, MUs can participate in the consensus process to obtain rewards as their mining

Without loss of generality, we consider that each MU $n$ has multiple IoT data processing tasks denoted by a set $\mathcal{M} = \{1, 2, ..., M\}$. Therefore, we denote $R_{nm}$ as the *m-th* data task of user n. Each miner is identified by a unique blockchain account. We also assume that the MEC server has sufficient computation resources with high-frequency CPU cores and storage capacity, and can provide computation services for a certain number of mobile miners in the blockchain network.

It is also assumed that time is slotted and at each time slot $t$, the miner $n$ generates a task $R_{nm}$. For each miner $n$, the data processing task at the timeslot $t$ can be formulated as a variable tuple $R_{nm}^t \triangleq (D_{1nm}^t, D_{0nm}^t, X_{nm}^t, \tau_{nm}^t) \in \mathcal{R}$. Here $D_{1nm}^t$ (in bits) denotes the data size of IoT data newly received from IoT devices by the miner $n$ at the timeslot $t$. $D_{0nm}^t$ (in bits) is the data size of IoT data in the buffer of the miner $n$. $X_{nm}^t$ (in CPU cycles/bit) denotes the total number of CPU cycles required to accomplish the computation for the task $R_{nm}^t$. Moreover, $\tau_{nm}^t$ (in seconds) reflects the maximum tolerable delay of task $R_{nm}^t$.

### B. Offloading For IoT Data Processing Tasks

In this paper, we consider a realistic blockchain scenario where the transactions from IoT devices are highly dynamic. That means at a certain time, the transaction volume can be very large (i.e., during the data collection of IoT devices) and at another time, this can be very small (i.e., idle mode of IoT devices). Based on that, at the timeslot $t$, if the size of data processing task is too large, the miners have to divide

contribution. However, due to the high computation requirement by the mining puzzle and resource constraints of mobile devices, cloud-based mining methodologies such as the leased hashing power technique [30] can be adopted which enables MUs to lease the required amount of hashing power from edge cloud servers. In this way, MU can purchase the necessary hash power from the edge cloud providers (e.g., AWS cloud services) for mining blocks. From the user point of view, the cloud-based mining process can eliminate any hassle of managing the infrastructure, i.e, installing extra mining hardware or requiring more computing resources on mobile devices. Furthermore, MUs can gain more economic benefits for their mobile blockchain apps through mining rewards, which will attract participation of more users and in return improve the robustness of the IoT blockchain network. Note that in our concerned blockchain network, MUs can leverage existing mining communication protocols such as Stratum [31] to achieve low-latency and secure mining process.

### D. Offloading Privacy

In the task offloading process, the MEC server is assumed to be curious about the personal information of miners. Motivated by [29], in this paper, we consider two privacy issues in the task offloading process, namely *user location privacy* and *usage pattern privacy* of the miner. By monitoring task offloading history of the miner, the MEC server can infer channel state records to obtain location information. Further, the MEC server can infer personal information of the miner by tracking mining data through the task offloading period, which may raise a concern of privacy threat based on the usage pattern of the miner. Thus, how to build an optimal task offloading policy that improves user privacy while guaranteeing low computation costs in the dynamic wireless environment is important to any MEC-based mobile blockchain applications. In the next subsection, we present the offloading problem formulation for the proposed model.

### IV. PROBLEM FORMULATION

In this section, we formulate the offloading problems for IoT data tasks and data processing tasks in details.

### A. Offloading Cost of IoT Data Processing Tasks

In this subsection, we derive the expression of performance metrics on network latency, energy consumption and user privacy under two offloading cost categories.

*1) Task offloading cost:* We model the case when the MU $n$ selects to offload its data processing tasks to the MEC server ($x_{nm} = 1$). Offloading latency and energy consumption problems are considered in our offloading formulation.

- *Offloading latency:* The computation latency is the key performance indicator for evaluating the quality of data processing task offloading. For the miner $n$, ($\forall n \in \mathcal{N}$), the whole computation latency to offload the data processing task $m$ to the MEC server at the timeslot $t$ can be decomposed into three parts, namely uploading, queuing and processing [11], which is given as

$$L_{nm}^{o,t} = L_{nm}^{u,t} + L_{nm}^{p,t} + L_{nm}^{q,t}, \qquad (1)$$

where $L_{nm}^{u,t}, L_{nm}^{p,t}, L_{nm}^{q,t}$ are uploading delay, task processing delay and queuing latency on MEC server, respectively.

- *Uploading delay:* We consider the delay when the user $n$ uploads the data processing task $m$ to the MEC server. We use $r_n$ to denote the upload transmission rate of the user $n$ in the wireless channel. For simplicity, we assume that the transmission rate is the same for all users $n$. Thus, we can express the required time to upload the data processing task $m$ to the MEC server as

$$L_{nm}^{u,t} = \frac{D_{0nm}^t + D_{1nm}^t}{r_n}. \qquad (2)$$

- *Queuing latency:* In our task offloading model, we are interested in the queuing latency caused by waiting in the task buffer for processing at the MEC server. Let $Q$ denote as the total number of CPU cycles in the data processing task buffer, the queuing latency can be given by

$$L_{nm}^{q,t} = \frac{Q}{f_{nm}}. \qquad (3)$$

- *Processing delay:* The time consumed by the miner $n$ to execute the data processing task $m$ on the MEC server at the timeslot $t$ can be formulated as [22]

$$L_{nm}^{p,t} = \frac{(D_{0nm}^t + D_{1nm}^t)X_{nm}^t}{f_{nm}}, \qquad (4)$$

where $f_{nm}$ is defined as the computational resource (in CPU cycles per second) allocated by the MEC server to accomplish the task $R_{nm}$. Note that this computation capacity is assumed to be large enough to serve all miners our blockchain system.

Once the execution of data processing task on MEC server finishes, the processed result is downloaded to the miner. Thus the downloading time latency can be given by $L_{nm}^d = \frac{G_{nm}}{S_{nm}}$ where $G_{nm}$ defines the size of executed data and $S_{nm}$ is the data rate of the MU $n$. However, the processed data from the MEC server is very small and the download rate is very high in general [12], and thus the latency (and energy cost) for the download process will be ignored in this paper.

- *Energy consumption:* The entire energy consumption for the task offloading can be formulated as the sum of energy costs for task uploading, task processing and MEC operation. Thus it is easy to express the total energy consumed by the computation offloading for the data processing task $m$ of the MU $n$ at the time slot $t$ as

$$E_{nm}^{o,t} = P^M L_{nm}^{u,t} + \gamma(f_{nm})^3 L_{nm}^{p,t} + P^C L_{nm}^{q,t}, \qquad (5)$$

where $\gamma$ is the energy consumption efficiency coefficient of the MEC server, $P^M$ is the transmit power of miner and $P^C$ denotes the baseline circuit power [13].

*2) Local computation cost:* When the MU $n$ decides to execute its task $m$ locally ($x_{nm} = 0$), it uses computation resource of the local device to process the mining puzzle. In this case, we consider the time cost and energy consumption for the local mining process.

We denote the local executing time per data bit as $t_{nm}^l$ (in sec/bit), then the total time cost consumed by the miner $n$ to complete the task $m$ at the timeslot $t$ can be expressed by

$$L_{nm}^{l,t} = (D_{0nm}^t + D_{1nm}^t)t_{nm}^l. \qquad (6)$$

Meanwhile, we use $e^l_{nm}$ (in J/bit) to denote the energy consumption per data bit of the MU $n$. Accordingly, the energy cost required for the MU $n$ to execute the data processing task at the timeslot $t$ is

$$E^{l,t}_{nm} = (D^t_{0nm} + D^t_{1nm})e^l_{nm}. \tag{7}$$

### B. User Privacy

In this subsection, we consider user privacy issues that have been largely ignored in previous studies in MEC-based mobile blockchain networks. Here, offloading privacy is defined as the preservation of the user information (including usage pattern and user location) during the mobile task offloading process against unintended usage or threats. Through the user's wireless transmission activity, an eavesdropper can monitor the service migrations among MEC and users and track physical movements of the user. An attacker can estimate the size of the sensing data newly generated and thus evaluate the usage pattern if the user offloads all the sensing data to the edge device under good radio channel state. Furthermore, the user also should prevent himself from being tracked by the eavesdropper or malicious MEC servers, and the offloading strategy needs to mislead the attackers to avoid any privacy harms on the user [32]. Thus, user information can be securely maintained which improves significantly user privacy, accordingly. Based on above discussion, here we consider two privacy issues for task offloading: usage pattern privacy and location privacy.

*1) Usage pattern privacy:* Naturally, if the wireless channel state is good, the miner tends to offload their data processing task with all blockchain transaction data ($D^t_{offload} = D^t_{1nm} + D^t_{0nm}$) to the MEC server to reduce processing time and energy cost on the mobile device. More specially, if the mobile user $n$ moves near to the access point, it is more likely to have $D^t_{offload} = D^t_{1nm}$ (due to $D^t_{0nm} = 0$). Obviously, the data usage pattern of the miner can be obtained by the MEC server through monitoring the data processing task $D^t_{offload}$. We denote $\zeta$ as the pre-defined good channel power gain state, similar to [29], the level of usage pattern privacy can be estimated by

$$P^{u,t}_{nm} = |D^t_{0nm} - x^t_{nm}(D^t_{0nm} + D^t_{1nm})|.\mathbb{I}(g^t_n \geq \zeta), \tag{8}$$

where $\mathbb{I}$ denotes the indicated function that equals 1 (or 0) if the statement is true (or false). The equation (8) means that the miner $n$ deliberately changes the amount of transaction data processed by local device under the good channel state, aiming to create a difference between the actual transaction data size $D^t_{0nm}$ and the offloading data size ($D^t_{0nm} + D^t_{1nm}$) to preserve usage pattern privacy.

*2) Location privacy:* Besides usage pattern privacy, the location privacy is another critical issue that should be considered to improve the performance of blockchain data processing task offloading. According to the system model formulation in the previous subsection, the MU will offload its data processing task to the MEC server if the wireless channel has a good transmission state ($g^t_{nm} = 1$). Otherwise, the data processing task will be executed locally or stored in the buffer for future process when the channel state is bad ($g^t_{nm} = 0$). In fact, we can encourage more mobile users to offload tasks to the MEC server when the channel state is good for efficient computation. However, since the wireless channel power gain is highly correlated to the distance between the miner and MEC server, this offloading strategy can expose location information of miner [32]. Applying the privacy metric in [29], we can formulate the location privacy level as

$$P^{l,t}_{nm} = \mathbb{I}[x^t_{nm}(D^t_{0nm} + D^t_{1nm})].\mathbb{I}(g^t_n < \zeta), \tag{9}$$

which means that the miner should keep offloading its data processing task to MEC server under poor transmission channel quality to preserve its location privacy. In summary, to protect privacy during task offloading, miners should mitigate the offloading rate under good channel quality while increasing the offloading rate under poor channel quality.

Combining the result (8) and (9), the total task offloading privacy level at the timeslot $t$ can be given as

$$P^t_{nm} = P^{u,t}_{nm} + \lambda P^{l,t}_{nm}, \tag{10}$$

where $\lambda$ scales the importance of the location privacy relative to the usage pattern privacy ($0 < \lambda < 1$).

### C. Reward of Blockchain Mining

As part of the blockchain network, MUs can join the consensus process to perform mining tasks for extra profits. In the blockchain network, MUs compete against each other to become the first one to solve the mining puzzle. Each MUs $n$ determines its mining service demand to purchase the hash power from the edge-cloud provider, denoted as $p_n$ (Hash/sec). Then, the relative hash power of the MU $n$ to the blockchain network can be defined as

$$\mu_n = \frac{p_n}{H}, \tag{11}$$

where $H$ represents the total hash power of the blockchain network which can be estimated through the compact status reports from miners [33]. It can be seen that when $\mu_n$ increases, the probability to achieve successful consensus increases, which will increase the mining reward for the MU $n$, accordingly.

In general, in order to mine successfully a block, two steps are needed including the mining step and the propagation step. In the mining step, the probability that the MU $n$ mines the block is proportional to its relative hash power $\mu_n$. The mined block is then propagated to the blockchain network. However, there is possibility that this MU $n$ propagates the mined block slower with other miners in the propagation step, which makes such a block likely to be discarded from the blockchain. This issue is called as orphaning [34], and this miner does not receive a reward for its mining. According to [35], the orphaning probability is approximated as $\mathcal{P}_{orphan} = 1 - e^{-\eta\phi(s_n)}$, where $\eta$ is a constant mean value and set $\eta = 1/600(sec)$ [36]. Further, $s_n$ denotes the number of transactions contained in the block that is mined by the miner $n$, and $\phi(s_n)$ is a function of block size [37], which represents the block propagation time. Clearly, with a larger block size $s_n$, the propagation time required to reach consensus of a block will be larger. As

a result, the probability of successful mining by miner $n$ can be expressed as

$$\mathcal{P} = \mu_n(1 - \mathcal{P}_{orphan}) = \mu_n e^{-\eta\phi(s_n)}. \quad (12)$$

We denote $\mathcal{R}$ as the reward of the first miner which achieves consensus, then the expected mining reward of the MU $n$ can be calculated as $R_n = \mathcal{R}\mathcal{P} = \mathcal{R}\mu_n e^{-\eta\phi(s_n)}$. Moreover, the process of solving the mining puzzle incurs an associated cost, i.e., a payment from the miner $n$ to the edge-cloud provider, denoted as $Y_n$ (token). Finally, based on the above mining formulation, the total mining reward of data processing tasks of the miner $n$ can be expressed as

$$R_n^{mining} = \mathcal{R}\mu_n e^{-\eta\phi(s_n)} - Y_n. \quad (13)$$

## V. OFFLOADING OPTIMIZATION WITH RL

In this section, we derive the optimization problem and then formulate a RL-based approach for the proposed scheme.

### A. Optimization Problem Formulation

In this section, we formulate the data processing task offloading and edge resource allocation as a joint optimization problem. The objective in this work is to maximize the offloading privacy while minimizing the sum cost of computation latency and energy consumption.

We first formulate the computation latency as the maximum of the local task processing time and the task execution time on MEC server, that can be expressed as

$$L_{nm}^t = \max\{L_{nm}^{o,t}, L_{nm}^{l,t}\}. \quad (14)$$

Meanwhile, the total energy consumption includes the local energy consumption $E_{nm}^{l,t}$ and energy cost for the taks offloading $E_{nm}^{o,t}$. Therefore, we can denote $C$ as the cost function that is the weighted sum of the time latency and energy consumption, as

$$C_{nm}^t = \sum_{n=1}^{N}\sum_{m=1}^{M}[\alpha_1(x_{nm}^t E_{nm}^{o,t} + (1 - x_{nm}^t)E_{nm}^{l,t}) + \alpha_2.L_{nm}^t], \quad (15)$$

where $\alpha_1, \alpha_2 \in (0, 1)$ denote the weight of energy consumption and task processing latency. Now we can formulate the optimization problem to jointly optimize the system privacy 10, system cost 15, and mining reward 13 under the constraint of maximum task mining latency and MEC computation capacity, which can be expressed as follows

$$P1: \quad \max_{\mathbf{x}} \sum_{n=1}^{N}\sum_{m=1}^{M}(P_{nm} + R_n^{mining} - C_{nm})$$

subject to

$$x_{nm} \in \{0, 1\}, \forall n \in \mathcal{N}, m \in \mathcal{M}, \quad (16a)$$

$$x_{nm}L_{nm}^o + (1 - x_{nm}L_{nm}^l) \leq \tau_{nm}, \forall n \in \mathcal{N}, m \in \mathcal{M}. \quad (16b)$$

Here, $\mathbf{x} = [x_{11}, x_{12}, ..., x_{1M}; ...; x_{N1}, x_{N2}, ..., x_{NM}]$ is the offloading decision vector. Here, the constraint 16a represents the binary offloading decision policy of the miner $n$ for the task

m, or offloading to the MEC server or processing locally at the mobile device. Further, the execution time to complete a data processing task should not exceed a maximum time latency value, which is expressed in the constraint 16b. It is worth noting that the optimization problem $(P1)$ is a mixed integer non-convex problem, which is NP-hard and difficult to derive an optimal solution. This is because the variable $x_{nm}$ is binary variables which make the problem $(P1)$ become a mixed integer programming problem that is in general non-convex and NP-hard. Further, the feasible set of problem $(P1)$ is also not convex. Also, the objective function is non-convex due to the product relationship between offloading decisions, the energy consumption factor and the offloading latency factor. To solve the offloading issues in the proposed task offloading, there are some challenges that need to be solved:

- The system states and rewards cannot be obtained in advance by step-wise control. Thus, the conventional optimization methods that only consider the current state may not be feasible.
- The joint optimization problem of offloading decision, offloading cost and mining profit with dynamicity of offloading data sizes, wireless channel states and hash power states is of high dimensionality and high complexity. It is hard to make joint efficient decisions by conventional approaches.
- As shown in the previous work [1], there are some regularities of user's preferences and network features. With these regularities, offloading strategies and mining strategies can be made in advance. Therefore, using deep reinforcement learning approach would be a viable solution to address the above challenges.

To overcome such challenges, we propose a dynamic offloading scheme using deep reinforcement learning (DRL). The advantages of our algorithm are twofold. First, it enables miners to obtain an optimal offloading action at each system state based on current blockchain transaction data and channel state without requiring prior knowledge of system dynamics. Second, as the DRL-based method is efficient in solving complex problems with large state space [18], it can achieve a better offloading performance to improve the quality of large-scale blockchain applications. Details of our design are presented in the next section.

### B. Reinforcement Learning Formulation

In our blockchain scenario, each miner acts an agent which interacts closely with the mobile blockchain environment at the timeslot $t$ to find a task offloading action $a$ for a state $s$ using a policy $\pi$ as shown in Fig. 3. This policy is defined as a mapping from the action to the state, i.e., $\pi(s) = a$. The main goal of the blockchain miner is to find an optimal policy $\pi$, aiming to maximize the total amount of award $r$ over the long run. To implement the RL-based algorithms, we first define the specific state, action and reward for the proposed task offloading. We formulate the main elements of a reinforcement learning approach for our blockchain offloading as follows.
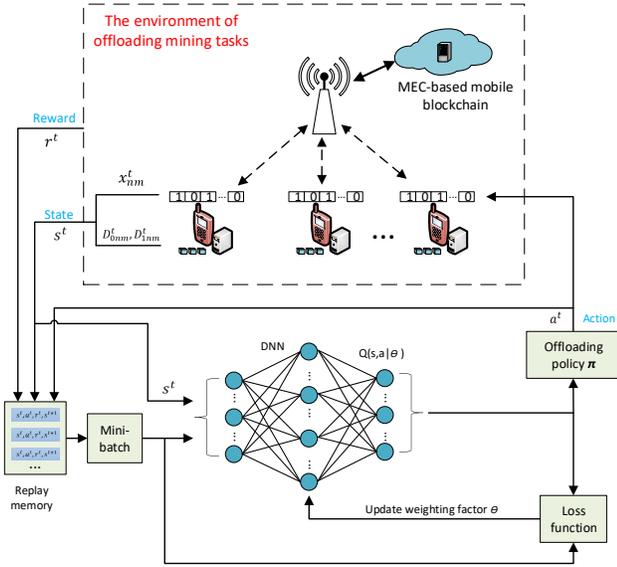
Fig. 3: RL-based offloading for mobile blockchain.

*1) State:* The system state is chosen as $s^t = \{D_{1nm}^t, D_{0nm}^t, g_n^t, p_n^t, Y_n^t\}$ where $D_{1nm}^t, D_{0nm}^t$ represent the new and buffered blockchain transaction data of the miner at the timeslot $t$, respectively. $g_n^t$ is the power gain state of wireless channel between the miner and MEC server as defined in the system model. Further, $p_n^t$ represents the hash power that the miner $n$ purchases from the edge-cloud provider, and $Y_n^t$ is the payment that the miner $n$ has to pay for the edge-cloud provider to perform mining.

*2) Action:* In our paper, power gain was pre-defined using the channel state transition probability, and the MEC computation capacity is assumed to be large enough to serve all miners in our blockchain system. For hash power and payment, each user determines its mining service demand to purchase the hash power from the edge-cloud provider before the data offloading starts, while payment is completed after the data offloading process in an automatic manner. As a result, in this paper we only need to consider the offloading policy of data tasks as the action space. It can be formulated as the offloading decision vector $\mathbf{x}^t = [x_{11}^t, x_{12}^t, ..., x_{1M}^t; ...; x_{N1}^t, x_{N2}^t, ..., x_{NM}^t]$. Therefore, the action vector can be expressed as $a^t = x_{nm}^t (\forall n \in \mathcal{N}, m \in \mathcal{M})$.

*3) Reward:* The objective of the RL agent is to find an optimal offloading decision action $a$ at each state $s$ with the aim of achieving the maximum mining reward $R(s,a)^{mining}$ and the highest privacy level $P(s,a)$ while minimizing the sum cost $C(s,a)$ of time and energy consumption in the task offloading. Specially, the reward function should be positively related to the objective function of the optimization problem (P1) in the previous section. Accordingly, we can formulate the immediate system reward $r^t(s,a)$ as

$$r^t(s,a) = P^t(s,a) + R_n^{mining} - C^t(s,a). \qquad (17)$$

In the next section, we propose task offloading schemes using Reinforcement Learning with two algorithms: RL-based task offloading (RLO) and deep RL-based task offloading (DRLO).

## VI. RL-BASED ALGORITHMS

### A. RL-based Task Offloading Algorithm

The principle of Reinforcement Learning (RL) can be described as a Markov Decision Process (MDP) [15]. In the RL model, an agent can make optimal actions by interacting with the environment without an explicit model of the system dynamics. In our blockchain scenario, we consider miners as agents to develop the RL scheme. At the beginning, the miner has no experience and information about the blockchain environment. Thus it needs to *explore* for every time epoch by taking some actions at each offloading state, i.e., the size of current blockchain transaction data. As long as the miner has some experiences from actual interactions with the environment, it will *exploit* the known information of states while keep exploration. In this paper, as a combination of Monte Carlo method and dynamic programing, a temporal-difference (TD) approach is employed to allow the agent to learn offloading policies without requiring the state transition probability which is difficult to acquire in realistic scenarios like in our dynamic mobile blockchain. In this subsection, we focus on developing a dynamic offloading scheme using a free-model RL. To this end, the state-action function can be updated using the experience tuple of agent $(s^t, a^t, r^t, s^{t+1})$ at each time step t as

$$Q(s^t, a^t) \leftarrow Q(s^t, a^t) + \alpha\sigma^t, \qquad (18)$$

which is called as Q-learning algorithm [15]. Here, $\sigma^t = r(s^t, a^t) + \gamma * \max Q(s^{t+1}, a^{t+1}) - Q(s^t, a^t)$ is the TD error which will be zero for the optimal Q-value. Further, $\alpha$ is the learning rate, $\gamma$ is the discount factor between (0,1). In line with the discussion, under the optimal policy $\pi^*$ which can be obtained from the maximum Q-value ($\pi^*(s) = argmaxQ^*(s,a)$), the Bellman optimality equation [15] for the state-action equation can be expressed as

$$Q^*(s^t, a^t) = \mathbb{E}_{s^{t+1} \sim E}[r(s^t, a^t) + \gamma * \max Q^*(s^{t+1}, a^{t+1})]. \qquad (19)$$

It is noting that the Q-learning algorithm is proved to converge with probability one over an infinite number of times [15] and achieves the optimal $Q^*$. Specially, we focus on finding the optimal policy $\pi^*$ which maximizes the reward $r(s,a)$ in equation 17. The details of task offloading implementation for mobile blockchain using Q-learning (RLO) are summarized in Algorithm 1.

Using Algorithm 1, the optimization problem (P1) can be solved to obtain the optimal offloading policy for miners. The RLO algorithm consists of two phases, the planning phase and the updating phase. The inputs are system states, i.e blockchain transaction data and channel states, and actions, i.e offloading tasks to MEC server or executing locally. The outputs are the resulting $Q^t(s,a)$ with maximum reward $r^*(s,a)$ which corresponds to the offloading policy $\pi^*(s,t)$ in each state. In the planning phase (lines 6-15), we use a $\epsilon$-greedy policy to balance the exploration and exploitation [15] for updating the Q function in equation 18. At each time epoch, the miner observes the blockchain state, selects an offloading action, and estimates the privacy value and system costs. After each action, the miner moves to the next step, updates the new state

---

**Algorithm 1** RL-based task offloading (RLO) algorithm for mobile blockchain networks

---
1: **Initialization:**
2: Initialize parameters: learning rate $\alpha$, discount factor $\gamma$, exploration rate $\epsilon \in (0,1)$
3: Initialize the action-value function Q with initial pair $(s,a)$, and create a Q-table with $Q(s^0, a^0)$
4: Set $t = 1$
5: **Procedure:**
6: **while** $t \leq T$ **do**
7:    /$* * *$ *Plan the task offloading in blockchain* $* * *$/
8:    Observe blockchain transaction $(D_1^t, D_0^t)$
9:    Estimate the channel gain state $g^t$
10:   Set $s^t = \{D_1^t, D_0^t, g^t\}$
11:   Select a random action $a^t$ with probability $\epsilon$, otherwise $a^t = argmaxQ(s^t, a, \theta)$
12:   Offload data processing task $x^t(D_1^t + D_0^t)$ to MEC server or execute $(1-x^t)(D_1^t + D_0^1)$ locally
13:   Calculate the system reward $r^t$ by equation 17
14:   Estimate the privacy level $P(s,a)^t$
15:   Estimate the system cost $C(s,a)^t$
16:   /$* * *$ *Update* $* * *$/
17:   Set $s^{t+1} = \{D_1^{t+1}, D_0^{t+1}, g^{t+1}\}$
18:   Update $Q(s^t, a^t)$ by equation 18
19:   $t \leftarrow t + 1$
20: **end while**

---

(lines 17-19) and iterates the offloading algorithm to obtain the optimal offloading policy.

Although the RLO algorithm can solve the problem (P1) by obtaining the optimum reward, there are still some remaining problems. The state and action values in the Q-learning method are stored in a two-dimensional Q table, but this method can become infeasible to solve complex problems with a much larger state-action space. This is because if we keep all Q-values in a table, the matrix $Q(s,a)$ can be very large, which makes the learning agents difficult to obtain sufficient samples to explore each state, leading to the failure of the learning algorithm. Moreover, the algorithm will converge very slow due to too many states that the agent has to process. To overcome such challenges, in the next subsection, we propose to use deep learning with Deep Neural Network (DNN) to approximate the Q-values instead of using the conventional Q-table.

### B. DRL-based Task Offloading Algorithm

In the DRL-based algorithm, a DNN is used to approximate the Q-values $Q(s^t, a, \theta)$ with weights $\theta$ as shown in Fig. 3. Further, to solve the instability of Q-network due to function approximation, the experience replay solution is employed in the training phase with the buffer $\mathcal{B}$ which stores experiences $e^t = (s^t, a^t, r^t, s^{t+1})$ at each time step $t$. Next, a random mini-batch of transitions from the replay memory is selected to train the Q-network. Here the Q-network is trained by iteratively updating the weights $\theta$ to minimize the loss function, which is written as

$$L^t(\theta^t) = \mathbb{E}[(r^t + \gamma * \max Q(s^{t+1}, a'|\theta') - Q(s^t, a^t|\theta^t))^2], \quad (20)$$

where the $\mathbb{E}[]$ denotes the expectation function. The detailed DQN-based task offloading algorithm for the proposed blockchain network is summarized in Algorithm 2.

---

**Algorithm 2** DRL-based task offloading (DRLO) algorithm for mobile blockchain networks

---
1: **Initialization:**
2: Set replay memory $\mathcal{D}$ with capacity $N$
3: Initialize the Q network with input pair $(s,a)$ and estimated action-value function $Q$ with random weight $\theta$; initialize the exploration probability $\epsilon \in (0,1)$
4: **for** $t = 1, 2, ...$ **do**
5:   /$* * *$ *Plan the task offloading in blockchain* $* * *$/
6:   Observe blockchain transaction $(D_1^t, D_0^t)$
7:   Estimate the channel gain state $g^t$
8:   Set $s^t = \{D_1^t, D_0^t, g^t\}$
9:   Select a random action $a^t$ with probability $\epsilon$, otherwise $a^t = argmaxQ(s^t, a, \theta)$
10:  Offload data processing task $x^t(D_1^t + D_0^t)$ to MEC server or execute $(1-x^t)(D_1^t + D_0^1)$ locally
11:  Observe the reward $r^t$ calculated via equation 17 and next state $s^{t+1}$
12:  Evaluate the achieved privacy $P(s,a)^t$, and system cost $C(s,a)^t$
13:  /$* * *$ *Update* $* * *$/
14:  Store the experience $(s^t, a^t, r^t, s^{t+1})$ into the memory $\mathcal{D}$
15:  Sample random mini-batch of state transitions $s^t, a^t, r^t, s^{t+1}$ from $\mathcal{D}$
16:  Calculate the Q-value by $y^k = r^k + \gamma * \max Q(s^{k+1}, a', \theta)$
17:  Perform a gradient descent step on $(y^k - Q(s^k, a^k, \theta))^2$ as loss function
18:  Train the Q-network with updated $\theta^t$
19: **end for**

---

The DRLO algorithm can achieve the efficient task offloading strategy in an iterative manner. As shown in Algorithm 2, the procedure generates a task offloading strategy for miners based on system states, e.g., blockchain transaction data, and observes the system reward at each time slot so that the offloading policy can be optimized (lines 6-12). Then the procedure updates the history experience tuple and train the Q-network (lines 14-18) with loss function minimization. This trial and error solution will avoid the requirement of prior information of offloading environment. Over the training time period, the trained neural network can characterize well the environment and therefore, the proposed offloading algorithm can dynamically adapt to the real mobile blockchain environment.

### VII. SIMULATION AND PERFORMANCE EVALUATION

In this section, we investigate the proposed offloading scheme by conducting both real experiments and simulations evaluate the efficiency of computation offloading.
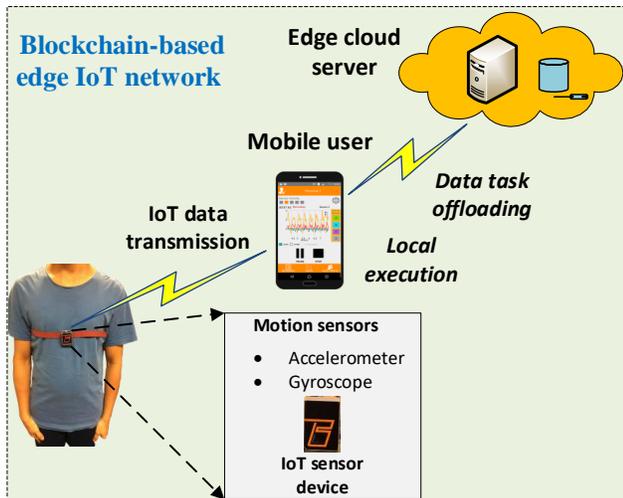
Fig. 4: Experiment setup for task offloading in blockchain.

## A. Implementation Settings

*1) Experiment settings:* We considered a task offloading framework for the MEC-based blockchain network as shown in Fig. 4. We deployed an Ethereum blockchain network supported by Amazon cloud where edge computing was deployed on the Lambda Edge[1] service for high transfer speeds and low-latency computation. A virtual machine Ubuntu 16.04 LTS, 2.4 GHz frequency was used as a computing unit for edge service. To test the proposed offloading scheme, we used a Sony mobile phone running on an Android OS version 8.0 platform with Qualcomm Snapdragon 845 processor, 4GB of RAM and 64GB of expandable storage, and a battery capacity of 2870mAh. The mobile devices connect with the edge cloud computing on the wireless network via Wi-Fi wireless communication with the standard IEEE 802.11g. A blockchain client was also installed on the device to transform the mobile phone into a miner node [3].

For IoT data generation, we used Biokin sensors [3] as IoT devices to collect health data which then needs to be computed for medical services. For the evaluation of local task computation on smartphones, we employed the Firebase Performance Monitoring[2] service to measure processing time and battery consumption. Meanwhile, for the evaluation of edge computing, we utilized the Kinesis Data Analytics service to monitor data streaming from mobile devices, and measure data computing latency and energy consumption. The evaluation needs IoT data and programming code for computation. We collected simultaneously data from sensors and stored them in separate files which are transmitted wirelessly to the mobile phone for computation. Details of hardware settings for our offloading can be found in our recent works [3], [38].

*2) Simulation settings:* In our simulation, a mobile blockchain network is considered with a MEC server with a varying number of mobile users (miners) and data processing tasks. The simulation is conducted over $T = 8,000$ timeslots, and each timeslot $t$ lasts 1 s. The computational capacity of

---

[1]https://aws.amazon.com/lambda/edge/
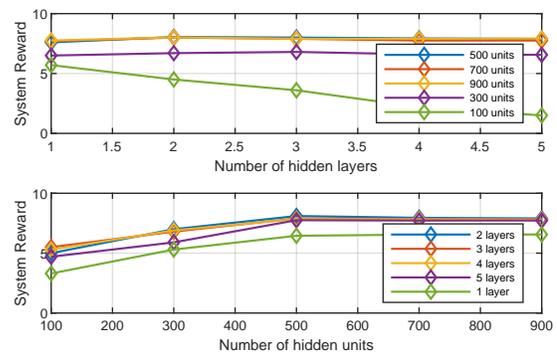[2]https://firebase.google.com/docs/perf-mon

---



Fig. 5: Impacts of different hyperparameters.

the MEC server $F$ is set to 10 GHz/sec. At each mobile user, we set the local computation time and computing energy consumption as $4.75 * 10^{-7}$ s/bit and $3.25 * 10^{-7}$ J/bit, respectively [9]. We assume that the size of data processing tasks generated from IoT devices is randomly distributed between 50 kB and 150 kB [37]. The local computation workload $X$ is set to 18000 CPU cycles/bit, and the delay threshold $\tau$ is assumed as 15s. The energy consumption efficiency coefficient and static circuit power of MEC server are set to $10^{-26}$ and 0.05W, respectively [22], and the channel gain factor $\sigma$ is set to 0.8. Further, based on [28], [29], the channel state transition probability is set to $Pr(g^{t+1} = 1|g^t = 1) = Pr(g^{t+1} = 0|g^t = 0) = 0.95$.

For the mining tasks, we configure the simulation parameters according to [37], [39], [40]. The size of block data to be mined is [5-10] kB. It is feasible in our blockchain scenario since data in the blockchain only includes metadata, while actual IoT data is stored in blockchain-basedd decentralized storage such as InterPlanetary File System (IPFS) [3]. The hash power $p_n$ that the MU $n$ purchases from the edge-cloud provider follows a uniform distribution on [20-100] MHash/s. The hash power of the blockchain network follows a uniform distribution on $10^3$-$10^5$ GHash/s. Furthermore, the miner that first solves the cryptographically hard puzzles and achieves successfully consensus is rewarded with $\mathcal{R} = 30$ tokens [37].

In the DQN-learning algorithm, we configure the parameters as follows. The discount factor $\gamma$ is set to 0.85; the learning rate equals 0.01, and the replay memory capacity and training batch size are set to $10^5$ and 128, respectively. We have employed gradient descent when approximating RL, and used a feedforward neural network to build our DNN model. Next, we employ ReLU as the activation function in the hidden layers, while the sigmoid activation function is utilized in the output layer to relax the offloading decision variables [38].

The architecture of the deep neural network significantly affects the performance of the deep reinforcement learning algorithms, and thus it requires a thoughtful design. We implemented simulations to select the parameter for our learning structure with the different DNN layers and varying hidden units. The experiment results in Fig. 5 assert that the number of hidden layers and neurons have a direct effect on the system throughput (e.g., system reward). Increasing the layers leads to better results, but at some point when the network is made deeper, the results start degrading. Their result shows that
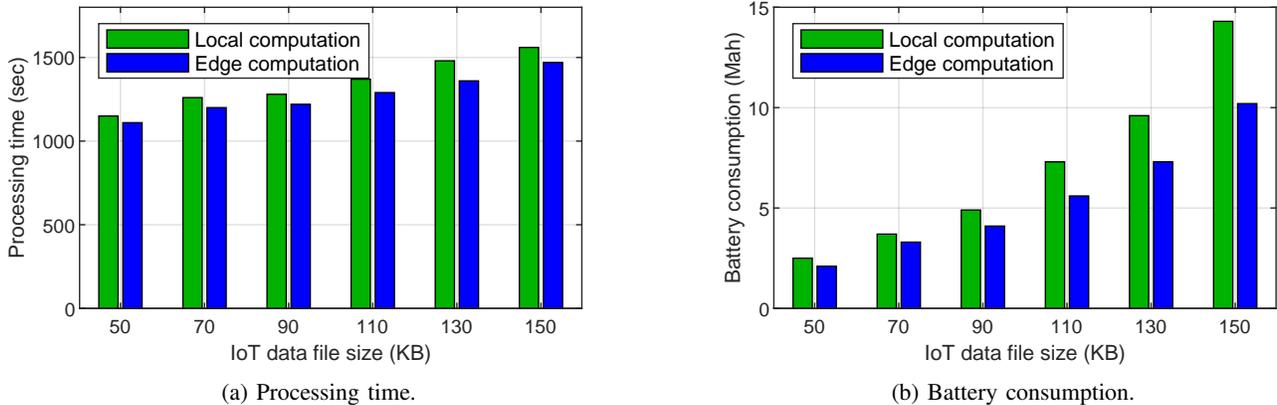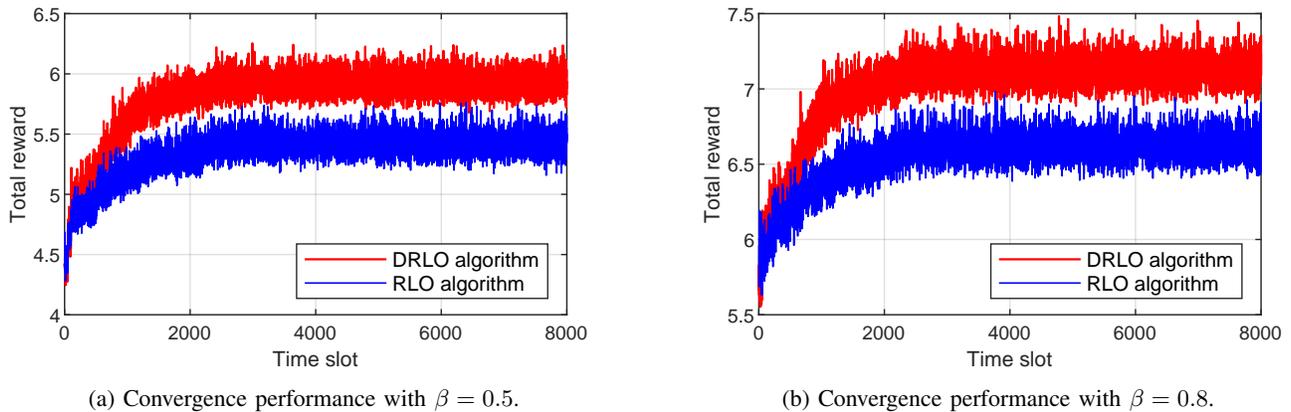
(a) Processing time.

(b) Battery consumption.

Fig. 6: Experimental results for local and edge computation.



(a) Convergence performance with $\beta = 0.5$.

(b) Convergence performance with $\beta = 0.8$.

Fig. 7: Convergence performance of offloading algorithms with $\beta = 0.5$ and with $\beta = 0.8$.

when using two hidden layers with 500 hidden units (including 300 units in first hidden layer and 200 units in second hidden layers), the model achieves the best performance. Hence, we choose a two-layer 500 DNN to build the offloading learning optimizer in this work.

To evaluate the task offloading performance for mobile blockchain, we focus on the following three metrics.

- The computation latency for offloading and local execution.
- The energy consumed by mining process on MEC server and local computation.
- The privacy level during the data processing task offloading process.

To highlight the advantage of the proposed offloading schemes in terms of latency and energy metrics, we compare our RLO and DRLO algorithms with two baseline schemes, i.e.,

- Non-offloading scheme (NO): All data processing tasks are executed at the local devices, (.i.e, setting offloading decision vector $x_{nm} = 0 (\forall n \in \mathcal{N}, m \in \mathcal{M})$).
- Edge offloading scheme (EO): All miners offload their data processing tasks to the MEC server (.i.e, setting offloading decision vector $x_{nm} = 1 (\forall n \in \mathcal{N}, m \in \mathcal{M})$).

Futher, to evaluate the privacy performance metric, we compare our algorithms with the constrained Markov decision process (CMDP)-based scheme [29] and RL-based design in [28].

Specially, we introduce a tradeoff factor $\beta \in [0, 1]$ for each miner to analyze the influence of computation laytency and energy consumption on the quality of mobile blockchain offloading. In this way, the weighted factors in equation 15 can be represented by $\alpha_1 = \beta$ and $\alpha_2 = 1 - \beta$. Therefore, the reward function equation 17 can be rewritten as
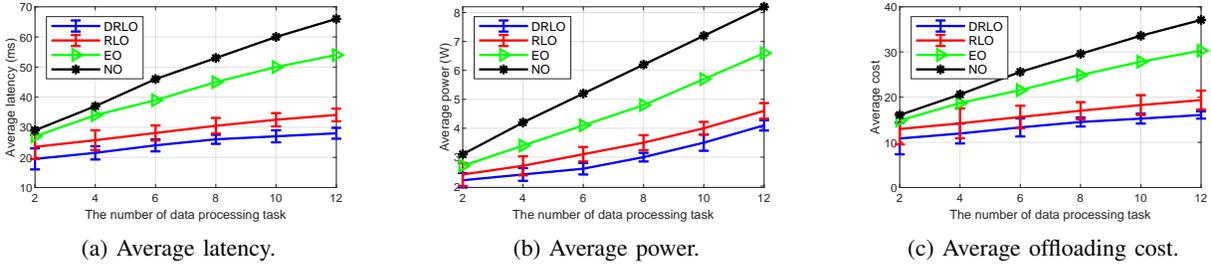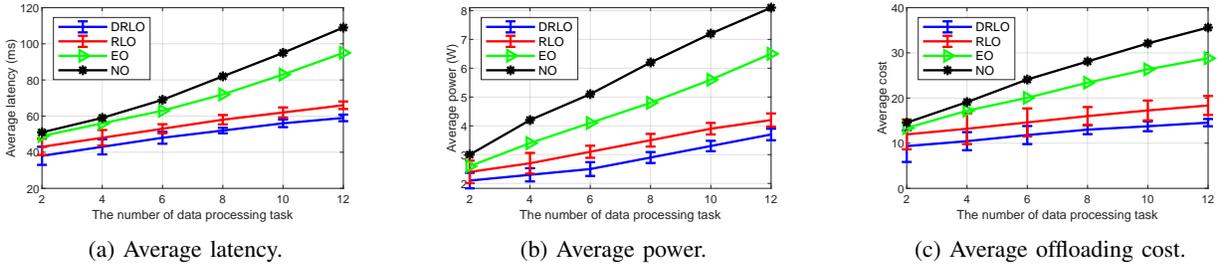
$$r^t(s, a) = P^t(s, a) - [\beta E^t(s, a) + (1 - \beta) L^t(s, a)], \quad (21)$$

where $E^t(s, a)$ and $L^t(s, a)$ are the sum offloading energy and latency cost, respectively.

### B. Experimental Evaluation

We considered two computation scenarios to prove the system efficiency: local execution and offloading to the edge computing service for computing IoT data. In the first scenario, the IoT data are executed locally on the mobile device. In the second scenario, the IoT data are offloaded for execution on the edge computing using our settings. A set of sensor data files with different sizes (50 kB-150 kB) was used to evaluate the proposed scheme. Each IoT data input is executed 10 times, and the average values are obtained.

We investigated the implementation results via two performance metrics: processing time and energy consumption as shown in Fig. 6. The processing time includes execution time, and offloading time with downloading time in the case of task offloading to the edge cloud server. In Fig. 6(a), the average processing time of local computation is higher than that of

(a) Average latency.　　　　　(b) Average power.　　　　　(c) Average offloading cost.

Fig. 8: Comparison results for single user scenario with $\beta = 0.5$.



(a) Average latency.　　　　　(b) Average power.　　　　　(c) Average offloading cost.

Fig. 9: Comparison results for single user scenario with $\beta = 0.8$.

edge computation for each IoT data file size. For example, the offloading scheme can save 36% time for completing computation of a 50 kB file and save up to 43% time for computing a 150 kB file, which shows advantages of the offloading scheme.

In the experimental result for battery consumption in Fig. 6(b), IoT data tasks consumes less energy when being executed with the offloading scheme because the resource-intensive computational tasks are offloaded to the edge-cloud servers. As an example, offloading the 50 kB file consumes less 15% energy than the case of local computation. Specially, the energy usage of the offloading scheme become more efficient when the data size increases. For instance, executing a 130 kB file and 150 kB file can save 17.8% and 19% energy, respectively when offloading the task to the edge-cloud server, compared to local execution.

### C. Simulation Evaluation

*1) Convergence performance:* We first evaluate the convergence performance of our data processing task offloading algorithms, namely RLO and DRLO with tradeoff factor settings $\beta = 0.5$ and $\beta = 0.8$ through the training process. Fig. 7 shows the learning curve obtained by the offloading policy using the proposed algorithms over 8,000 timeslots. The results show that the total system reward is very small at the beginning of the learning process. However, as the number of timeslot increases, the total reward increases rapidly and become stable after about 2500 timeslots, which validates the convergence performance of the proposed DQN-learning scheme. Further, it can be seen that the DRLO algorithm can achieve a better long-term reward in both cases, compared to the RLO scheme. For instance, in case of $\beta = 0.5$ in Fig. 7(a), the DRLO-based learning strategy converges to 5.9, which is approximately 5.2% higher than that of the RLO-based learning scheme, which converges to 5.4 after about 2500 timeslots. Moreover, for a larger tradeoff value with

$\beta = 0.8$ in Fig. 7(b), the DRLO scheme converges to about 7.1, which is roughly 9% higher than that of the RLO scheme. Overall, the DRLO scheme achieves better performance for all cases. The key reason behind that is the RLO with Q-learning always chooses system actions in a greedy manner. On the contrary, DRLO with deep Q-learning applies trial-and-error search to balance exploration and exploitation. In addition, the DRLO scheme with a DNN as an approximator outperforms in obtaining more valuable features of certain parts of states, by using the iterative target network. This also prevents the instabilities to propagate quickly and reduces the risk of divergence. As a result, the DRLO scheme can learn better from the state space for making better decisions which leads to an improvement of system reward. Next, we verify the proposed schemes via single user and multiple user scenarios. The simulation results are averaged from 50 runs of numerical simulations.

*2) Task offloading in single user scenario:* In the single user model, we consider a blockchain network with a single miner $N = 1$ and change the number of tasks at each miner ($M$=2∼12). Note that the size of each task is the same for all miners. The simulation results of the offloading performance in case of $\beta = 0.5$ and $\beta = 0.8$ are shown in Fig. 8 and Fig. 9. In Fig. 8 with $\beta = 0.5$, it is observed that when the number of task increases, the cost of four approaches increases due to the growing size of transaction data on the blockchain network. For example, in the EO strategy, the average offloading cost of all miners increases from 15 at $M = 2$ to 30 at $M = 12$, and that of the RLO scheme also increases to 9 at $M = 12$. Further, the NO algorithm always incurs a higher energy cost than the other schemes and has the largest increasing rate of 125%, reaching 36 at $M = 12$. The reason behind this observation is that when the number of data processing task grows, the computational capacity of the mobile miner becomes less sufficient to provide mining services for computing all tasks. Thus newly generated blockchain transactions have to wait to be processed in the buffer of local devices, which thus

TABLE II: **Comparison results for multi-user scenario with $\beta = 0.5$.**

| Schemes | Average power (W) | | Average latency(ms) | | Average cost | |
|---------|-------|--------|-------|--------|-------|--------|
| | N=5 | N=10 | N=5 | N=10 | N=5 | N=10 |
| DRLO | **3.5** | **6.9** | **32.6** | **55.8** | **19.8** | **34.8** |
| RLO | 4.2 | 8.6 | 35.3 | 63.2 | 21.9 | 40.2 |
| EO | 6.2 | 14.3 | 53.2 | 94.6 | 32.8 | 53.5 |
| NO | 4.5 | 9.1 | 39.7 | 75.4 | 24.6 | 46.8 |

TABLE III: **Comparison results for multi-user scenario with $\beta = 0.8$.**

| Schemes | Average power (W) | | Average latency(ms) | | Average cost | |
|---------|-------|--------|-------|--------|-------|--------|
| | N=5 | N=10 | N=5 | N=10 | N=5 | N=10 |
| DRLO | **3.3** | **6.5** | **55.9** | **90.7** | **13.8** | **23.3** |
| RLO | 3.9 | 8.3 | 59.1 | 98.5 | 14.9 | 26.3 |
| EO | 6.0 | 13.9 | 82.2 | 163.8 | 21.2 | 35.6 |
| NO | 4.2 | 8.5 | 66.3 | 116.7 | 16.6 | 30.1 |

increases the mining latency. Further, the higher cost comes from a higher computation delay and a power consumption. Specially, the DRLO scheme achieves the best performance with minimum sum cost for all data processing tasks.

Fig. 9 shows the simulation results with a larger tradeoff factor $\beta = 0.8$. According to the formulated reward function (18), a larger tradeoff factor will give more penalty to energy consumption. In this case, the gap between the DRLO scheme and the other baselines is larger than that in the case of $\beta = 0.5$, which is caused by the increased offloading delay and lower power cost. Particularly, the DRLO algorithm exhibits the best performance again among four schemes with minimum offloading efficiency index. For instance, when the miner has 12 data processing tasks, the offloading cost averaged over 50 simulations of the DRLO scheme is 18.7%, 57%, and 65% lower than those of the RLO, EO, and NO schemes, respectively. Such results demonstrate the efficiency of the DRLO-based scheme in the data task offloading.

*3) Task offloading in multi-user scenario:* Next, we analyze the task offloading performance for the blockchain network with multiple users. We consider two cases, $N = 5$ and $N = 10$. We also assume that there are 2 data processing tasks at each miner, and the size of each task is the same for all miners. The comparison results are shown in Table II and Table III for tradeoff factor $\beta = 0.5$ and $\beta = 0.8$, respectively. Under the scenario with $\beta = 0.5$ in Table II, the DRLO scheme exhibits the lowest average power consumption, computation latency in both cases $N = 5$ and $N = 10$, and thus achieves the minimum total offloading cost, followed by the RLO scheme with a small gap. Another observation is that among the four schemes, the highest offloading cost comes from the EO-based scheme, instead of the NO-based scheme. This is because the more miners offload data processing tasks to MEC server, the more computing capacity is required to provide enough computation resources for running the data processing tasks of all blockchain users. It is also noteworthy that the computation capacity of the MEC server is only sufficient to provide resources for a certain number of miners and high mining demands from multiple users can result in a significant increase in network latency and system cost.

Using a higher tradeoff value $\beta = 0.8$, we present comparison results in Table III. Similar to the single user scenario, the larger tradeoff factor gives more penalty to energy consumption, and therefore, the energy consumed by all miners is lower than that of the case of $\beta = 0.5$, while the offloading delay becomes larger for both cases of $N = 5$
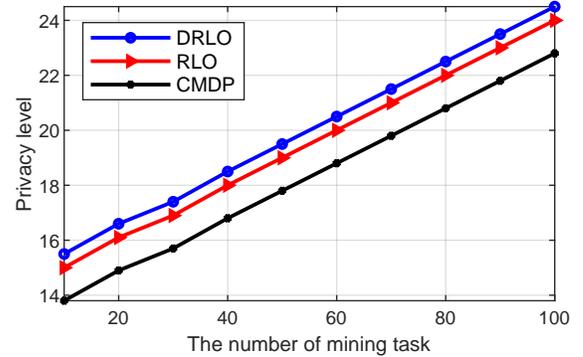


Fig. 10: The achieved privacy level.

and $N = 10$. Based on such observations, we can minimize energy consumption with respect to task offloading latency by adjusting the tradeoff factor for a better offloading efficiency. It is also worth mentioning that the DRLO-based scheme still exhibits the best offloading performance with minimum power usage, mining delay and resulting offloading cost. For instance, in the case of $N = 10$, the offloading cost of the DRLO-based strategy is 12.7%, 52.5%, and 30.4% lower than those of the RLO, EO, and NO schemes, respectively. The numerical results clearly show that the DRLO algorithm outperforms the other baseline schemes in improving offloading cost efficiency in multi-user scenarios.

*4) Privacy performance:* Finally, we analyse the performance of the proposed offloading design in terms of the privacy metric as shown in Fig. 10. We compare the privacy performance of our DRLO-based scheme with the RL-based scheme [29] and CMDP-based scheme [28]. It can be observed from Fig. 10 that the privacy level of the blockchain miner increases when the amount of transaction data increases from 10 kB to 100 kB in each time slot for all offloading schemes. However, the proposed DRLO method can achieve the best privacy performance, compared to the other benchmarks. For instance, as the mobile user mines 10 kB blockchain transactions, the proposed DRLO scheme achieves 5.2% and 12.7% higher privacy levels compared with the RL-based and CMDP-based schemes, respectively. Further, in the case of mining 100 kB blockchain transactions, the privacy level of DRLO scheme still shows the best performance, with 5.5% higher than the RL-based strategy and 13.4% higher than the CMDP-based strategy.

*D. Discussion and Future Works*

In this paper, we proposed a DRL-based dynamic task offloading scheme for a MEC blockchain network. Our focus was on formulating task offloading and user privacy preservation as a joint optimization problem. We then developed a RL-based algorithm to solve the proposed optimization problem with a simplified blockchain model. To this end, we conducted both experiments and numerical simulations to verify the effectiveness of the proposed offloading algorithm in terms of various performance metrics and compare with other offloading schemes. The simulation results with our proposed schemes showed the performance of computation offloading

for mobile blockchain networks can be significantly improved in terms of enhanced privacy level and reduced system costs, compared with other blockchain offloading schemes.

As future works, we will extend the proposed DRL scheme by taking action space variation into consideration. In this paper, we assume that the action space is stable during the learning process, but in certain blockchain scenarios, the action set can vary over time due to the dynamicity of the MEC-based blockchain offloading system, i.e., dynamic user demands, dynamic MEC resource usage, dynamic wireless channel power usage. In this regard, developing adative and online deep RL schemes would be helpful to cope with the dynamic environment. Another direction is to consider the dynamics of the network with the variable user set. In this case, the number of users in our blockchain network might vary with time, i.e., the current users can exit the network and new users can join the network. To cope with the dynamics of the network, one solution can be applied as follows. We can set a training value so that the remaining user states in the DNN network are regarded as empty states if the number of users is smaller than that preset training value. For instance, assuming that the preset training value is $X_p$ and the number of current users is $X_u (X_p > X_u)$, we can set the remaining $X_p - X_u$ states to be 0 and drop the remaining $X_p - X_u$ actions in the output of the DRLO algorithm. On the other hand, if $X_p < X_u$, we portion the collected user states into several parts and input these parts to the actor network in turn. This solution enables flexible adjustment on the learning algorithm according to the dynamics of the network [37].

## VIII. CONCLUSIONS

In this paper, we have proposed reinforcement learning-based task offloading algorithms for multi-users to obtain the optimal offloading policy in a dynamic blockchain network with MEC. We have formulated the task offloading and privacy preservation as a joint optimization problem. A RL-based scheme using the Q-network algorithm is employed to learn efficiently the offloading policy such that the total system cost combining computation latency and energy consumption is minimized while guaranteeing the best user privacy and mining reward performance. To break the curse of high dimensionality in state space, we have then developed a DRL-based approach using a deep Q-network algorithm. The offloading performances in terms of energy consumption, computation latency, and user privacy are analyzed under various conditions for both single user and multiple user offloading scenarios via numerical simulations. The experimental results have clearly showed that the proposed DRL offloading scheme is superior to the other baseline methods with a reduced energy consumption, computation latency with much lower offloading costs and improved privacy level.

## REFERENCES

[1] T. M. Fernandez-Carames and P. Fraga-Lamas, "A Review on the Use of Blockchain for the Internet of Things," *IEEE Access*, vol. 6, pp. 32 979–33 001, 2018.

[2] D. C. Nguyen, P. N. Pathirana, M. Ding, and A. Seneviratne, "Integration of Blockchain and Cloud of Things: Architecture, Applications and Challenges," *arXiv:1908.09058*, Aug. 2019.

[3] K. Suankaewmanee, D. T. Hoang, D. Niyato, S. Sawadsitang, P. Wang, and Z. Han, "Performance Analysis and Application of Mobile Blockchain," in *International Conference on Computing, Networking and Communications (ICNC)*, Maui, HI, Mar. 2018, pp. 642–646.

[4] D. C. Nguyen, P. N. Pathirana, M. Ding, and A. Seneviratne, "Blockchain for Secure EHRs Sharing of Mobile Cloud Based E-Health Systems," *IEEE Access*, vol. 7, pp. 66 792–66 806, 2019.

[5] R. Pass and E. Shi, "FruitChains: A Fair Blockchain," in *Proceedings of the ACM Symposium on Principles of Distributed Computing*, ser. PODC '17, Washington, DC, USA, Jul. 2017, pp. 315–324.

[6] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge Computing: Vision and Challenges," *IEEE Internet of Things Journal*, vol. 3, no. 5, pp. 637–646, Oct. 2016.

[7] P. Mach and Z. Becvar, "Mobile Edge Computing: A Survey on Architecture and Computation Offloading," *IEEE Communications Surveys & Tutorials*, vol. 19, no. 3, pp. 1628–1656, 2017.

[8] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A Survey on Mobile Edge Computing: The Communication Perspective," *IEEE Communications Surveys & Tutorials*, vol. 19, no. 4, pp. 2322–2358, 2017.

[9] H. Guo, J. Liu, J. Zhang, W. Sun, and N. Kato, "Mobile-Edge Computation Offloading for Ultradense IoT Networks," *IEEE Internet of Things Journal*, vol. 5, no. 6, pp. 4977–4988, Dec. 2018.

[10] M. Chen and Y. Hao, "Task Offloading for Mobile Edge Computing in Software Defined Ultra-Dense Network," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 3, pp. 587–597, Mar. 2018.

[11] S. Yu, X. Wang, and R. Langar, "Computation offloading for mobile edge computing: A deep learning approach," in *28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, Montreal, QC, Oct. 2017, pp. 1–6.

[12] X. He, H. Xing, Y. Chen, and A. Nallanathan, "Energy-Efficient Mobile-Edge Computation Offloading for Applications with Shared Data," in *IEEE Global Communications Conference (GLOBECOM)*, Abu Dhabi, United Arab Emirates, Dec. 2018, pp. 1–6.

[13] S. Li, N. Zhang, S. Lin, L. Kong, A. Katangur, M. K. Khan, M. Ni, and G. Zhu, "Joint Admission Control and Resource Allocation in Edge Computing for Internet of Things," *IEEE Network*, vol. 32, no. 1, pp. 72–79, Jan. 2018.

[14] I.-S. Comsa, S. Zhang, M. E. Aydin, P. Kuonen, Y. Lu, R. Trestian, and G. Ghinea, "Towards 5G: A Reinforcement Learning-Based Scheduling Solution for Data Traffic Management," *IEEE Transactions on Network and Service Management*, vol. 15, no. 4, pp. 1661–1675, Dec. 2018.

[15] R. S. Sutton and A. G. Barto, *Reinforcement learning: an introduction*, second edition ed., ser. Adaptive computation and machine learning series. Cambridge, Massachusetts: The MIT Press, 2018.

[16] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep Reinforcement Learning: A Brief Survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, Nov. 2017.

[17] 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.

[18] 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 Internet of Things Journal*, vol. 6, no. 3, pp. 4005–4018, Jun. 2019.

[19] J. Li, H. Gao, T. Lv, and Y. Lu, "Deep reinforcement learning based computation offloading and resource allocation for MEC," in *IEEE Wireless Communications and Networking Conference (WCNC)*, Barcelona, Apr. 2018, pp. 1–6.

[20] M. Min, L. Xiao, Y. Chen, P. Cheng, D. Wu, and W. Zhuang, "Learning-Based Computation Offloading for IoT Devices With Energy Harvesting," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 2, pp. 1930–1941, Feb. 2019.

[21] Z. Xiong, Y. Zhang, D. Niyato, P. Wang, and Z. Han, "When Mobile Blockchain Meets Edge Computing," *IEEE Communications Magazine*, vol. 56, no. 8, pp. 33–39, Aug. 2018.

[22] M. Liu, F. R. Yu, Y. Teng, V. C. M. Leung, and M. Song, "Computation Offloading and Content Caching in Wireless Blockchain Networks With Mobile Edge Computing," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 11, pp. 11 008–11 021, Nov. 2018.

[23] Y. Wu, X. Chen, J. Shi, K. Ni, L. Qian, L. Huang, and K. Zhang, "Optimal Computational Power Allocation in Multi-Access Mobile Edge Computing for Blockchain," *Sensors*, vol. 18, no. 10, p. 3472, Oct. 2018.
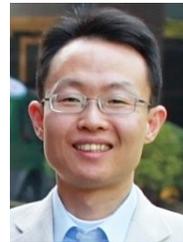
[24] M. Liu, F. R. Yu, Y. Teng, V. C. M. Leung, and M. Song, "Distributed Resource Allocation in Blockchain-Based Video Streaming Systems With Mobile Edge Computing," *IEEE Transactions on Wireless Communications*, vol. 18, no. 1, pp. 695–708, Jan. 2019.

[25] J. Zhang, B. Chen, Y. Zhao, X. Cheng, and F. Hu, "Data Security and Privacy-Preserving in Edge Computing Paradigm: Survey and Open Issues," *IEEE Access*, vol. 6, pp. 18 209–18 237, 2018.

[26] T. Salman, M. Zolanvari, A. Erbad, R. Jain, and M. Samaka, "Security Services Using Blockchains: A State of the Art Survey," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 1, pp. 858–880, 2019.

[27] D. K. Tosh, S. Shetty, X. Liang, C. A. Kamhoua, K. A. Kwiat, and L. Njilla, "Security Implications of Blockchain Cloud with Analysis of Block Withholding Attack," in *17th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID)*, Madrid, Spain, May 2017, pp. 458–467.

[28] X. He, J. Liu, R. Jin, and H. Dai, "Privacy-Aware Offloading in Mobile-Edge Computing," in *IEEE Global Communications Conference*, Singapore, Dec. 2017, pp. 1–6.

[29] M. Min, X. Wan, L. Xiao, Y. Chen, M. Xia, D. Wu, and H. Dai, "Learning-Based Privacy-Aware Offloading for Healthcare IoT With Energy Harvesting," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 4307–4316, Jun. 2019.

[30] D. K. Tosh, S. Shetty, X. Liang, C. A. Kamhoua, K. A. Kwiat, and L. Njilla, "Security Implications of Blockchain Cloud with Analysis of Block Withholding Attack," in *17th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID)*, Madrid, Spain, May 2017, pp. 458–467.

[31] R. Recabarren and B. Carbunar, "Hardening Stratum, the Bitcoin Pool Mining Protocol," *arXiv:1703.06545*, Mar. 2017.

[32] T. He, E. N. Ciftcioglu, S. Wang, and K. S. Chan, "Location Privacy in Mobile Edge Clouds: A Chaff-Based Approach," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 11, pp. 2625–2636, Nov. 2017.

[33] A. P. Ozisik, G. Bissias, and B. Levine, "Estimation of Miner Hash Rates and Consensus on Blockchains (draft)," *arXiv:1707.00082*, Jun. 2017.

[34] D. C. Nguyen, P. N. Pathirana, M. Ding, and A. Seneviratne, "Blockchain for 5G and beyond networks: A state of the art survey," *Journal of Network and Computer Applications*, vol. 166, p. 102693, Sep. 2020.

[35] Z. Li, Z. Yang, S. Xie, W. Chen, and K. Liu, "Credit-Based Payments for Fast Computing Resource Trading in Edge-Assisted Internet of Things," *IEEE Internet of Things Journal*, vol. 6, no. 4, pp. 6606–6617, Aug. 2019.

[36] N. Houy, "The Bitcoin Mining Game," *Ledger*, vol. 1, pp. 53–68, Dec. 2016.

[37] X. Qiu, L. Liu, W. Chen, Z. Hong, and Z. Zheng, "Online Deep Reinforcement Learning for Computation Offloading in Blockchain-Empowered Mobile Edge Computing," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 8, pp. 8050–8062, Aug. 2019.

[38] D. C. Nguyen, P. N. Pathirana, M. Ding, and A. Seneviratne, "Secure Computation Offloading in Blockchain based IoT Networks with Deep Reinforcement Learning," *arXiv:1908.07466*, Aug. 2019.

[39] S. Biswas, K. Sharif, F. Li, B. Nour, and Y. Wang, "A Scalable Blockchain Framework for Secure Transactions in IoT," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 4650–4659, Jun. 2019.

[40] Q. Luo, C. Li, T. H. Luan, and W. Shi, "Collaborative Data Scheduling for Vehicular Edge Computing via Deep Reinforcement Learning," *IEEE Internet of Things Journal*, pp. 1–1, 2020.

**Pubudu N. Pathirana** (Senior Member, IEEE) was born in 1970 in Matara, Sri Lanka, and was educated at Royal College Colombo. He received the B.E. degree (first class honors) in electrical engineering and the B.Sc. degree in mathematics in 1996, and the Ph.D. degree in electrical engineering in 2000 from the University of Western Australia, all sponsored by the government of Australia on EMSS and IPRS scholarships, respectively. He was a Postdoctoral Research Fellow at Oxford University, Oxford, a Research Fellow at the School of Electrical Engineering and Telecommunications, University of New South Wales, Sydney, Australia, and a Consultant to the Defence Science and Technology Organization (DSTO), Australia, in 2002. He was a visiting professor at Yale University in 2009. Currently, he is a full Professor and the Director of Networked Sensing and Control group at the School of Engineering, Deakin University, Geelong, Australia and his current research interests include Bio-Medical assistive device design, human motion capture, mobile/wireless networks, rehabilitation robotics and radar array signal processing.

**Ming Ding** (Senior Member, IEEE) received the B.S. and M.S. degrees (Hons.) in electronics engineering and the Ph.D. degree in signal and information processing from Shanghai Jiao Tong University (SJTU), Shanghai, China, in 2004, 2007, and 2011, respectively. From April 2007 to September 2014, he worked as a Researcher/Senior Researcher/Principal Researcher at the Sharp Laboratories of China, Shanghai. He also served as the Algorithm Design Director and the Programming Director for a system-level simulator of future telecommunication networks in Sharp Laboratories of China for more than seven years. He is currently a Senior Research Scientist with the CSIRO Data61, Sydney, NSW, Australia. His research interests include information technology, data privacy and security, machine Learning and AI. He has authored over 100 articles in IEEE journals and conferences, all in recognized venues, and around 20 3GPP standardization contributions, and a Springer book Multi-Point Cooperative Communication Systems: Theory and Applications. He holds 21 U.S. patents and co-invented another more than 100 patents on 4G/5G technologies in CN, JP, KR, EU. He is an Editor of the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS and the IEEE Wireless Communications Letters. Besides, he is or has been a Guest Editor/CoChair/Co-Tutor/TPC Member of several IEEE top-tier journals/conferences, such as the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS, IEEE Communications Magazine, and the IEEE GLOBECOM Workshops. He was the Lead Speaker of the industrial presentation on unmanned aerial vehicles in IEEE GLOBECOM 2017, which was awarded as the Most Attended Industry Program in the conference. He was awarded as the Exemplary Reviewer of the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS in 2017.

**Aruna Seneviratne** (Senior Member, IEEE) is currently a Foundation Professor of telecommunications with the University of New South Wales, Australia, where he holds the Mahanakorn Chair of telecommunications. He has also worked at a number of other Universities in Australia, U.K., and France, and industrial organizations, including Muirhead, Standard Telecommunication Labs, Avaya Labs, and Telecom Australia (Telstra). In addition, he has held visiting appointments at INRIA, France. His current research interests are in physical analytics: technologies that enable applications to interact intelligently and securely with their environment in real time. Most recently, his team has been working on using these technologies in behavioral biometrics, optimizing the performance of wearables, and the IoT system verification. He has been awarded a number of fellowships, including one at British Telecom and one at Telecom Australia Research Labs.

**Dinh C. Nguyen** (Graduate Student Member, IEEE) is currently pursuing the Ph.D. degree at the School of Engineering, Deakin University, Victoria, Australia. His research interests focus on blockchain, deep reinforcement learning, mobile edge/cloud computing, network security and privacy. He is currently working on blockchain and reinforcement learning for Internet of Things and 5G networks. He has been a recipient of the prestigious Data61 PhD scholarship, CSIRO, Australia.