

When User-Centric Network Meets Mobile Edge Computing: Challenges and Optimization Research Proposal

Langtian Qin, Hancheng Lu, and Feng Wu

August 17, 2022

1 System Model

1.1 Network Model

In this paper, we assume there are M single antenna users and N single antenna APs in UCMEC system, the user set and AP set are represented by $\mathcal{M} = \{1, 2, \dots, m, \dots, M\}$, $\mathcal{N} = \{1, 2, \dots, n, \dots, N\}$, respectively. For user m , $m \in \mathcal{M}$, it has a computation intensive task L_m . All APS are connected to a CPU integrated with MEC server through a reliable and unlimited bandwidth wired backhaul. Among them, the AP is only responsible for receiving and transmitting signals, while the CPU performs active user detection (AUD), channel estimation (CE) and calculation processing of user offloading tasks. Each user can choose to use the local limited computing resource for task calculation, or offload the tasks to the CPU with more powerful computing resource for calculation. We assume that the local computing resource of user m is C_m^l (in CPU-cycle frequency), while the computing resource of MEC server is C_e .

1.2 Channel and Communication Model

We assume that all communications are carried out in the same frequency band, and the UCMEC system operates in a time division multiplexed (TDD) mode. All users can send orthogonal pilot sequences at the same time. Assume $g_{m,n}$ is the channel between user m and AP n , which remains unchanged in the coherence time and satisfies the following formula

$$g_{m,n} = \sqrt{\beta_{m,n}} h_{m,n}, \quad (1)$$

where $\beta_{m,n}$ is the scalar coefficients of channel shadow fading (large scale fading) between user m and AP n , $h_{m,n}$ is the coefficient of small scale channel fading between user m and AP n , which satisfies $h_{m,n} \sim \mathcal{CN}(0, 1)$. Similar to [1], $\beta_{m,n}$ can be expressed by

$$\beta_{m,n} = 10^{\frac{\mu_{m,n}}{10}} 10^{\frac{\sigma_s \theta_{m,n}}{10}}, \quad (2)$$

where $\mu_{m,n}$ is the path loss (in dB) between user m and AP n , $10^{\frac{\sigma_s \theta_{m,n}}{10}}$ is the shadow fading with standard deviation σ_s between user m and AP n , $\theta_{m,n}$ is the shadow fading coefficient between user m and AP n . We use the following three slope model[2] to describe the path loss

$$\begin{cases} -L - 35\log_{10}(d_{m,n}), & d_{m,n} > d_1 \\ -L - 10\log_{10}(d_1^{1.5}d_{m,n}^2), & d_0 < d_{m,n} \leq d_1 \\ -L - 10\log_{10}(d_1^{1.5}d_0^2), & d_{m,n} < d_0, \end{cases} \quad (3)$$

where $d_{m,n}$ is the distance (in km) between user m and AP n , d_0 and d_1 are the breakpoints of three slope model. L is given by

$$L = 46.3 + 33.9\log_{10}(f) - 13.82\log_{10}(h_a) - [1.11\log_{10}(f) - 0.7]h_u + 1.56\log_{10}(f) - 0.8, \quad (4)$$

where f is the carrier frequency (in MHz), h_a and h_u indicates the antenna heights of the AP and the user, respectively. Since there are obstacles between the AP and the user in the real environment, through the correction of shadow fading[3], $\theta_{m,n}$ consists of two parts

$$\theta_{m,n} = \sqrt{\delta}a_n + \sqrt{1-\delta}b_m, m = 1, 2, \dots, M, n = 1, 2, \dots, N, \quad (5)$$

where a_n 和 b_m obey Gaussian distribution $\mathcal{N}(0, 1)$, $0 \leq \delta \leq 1$ is a parameter.

Similar to [4], We assume that at the beginning of each coherence time t_c , user m need to randomly choose a pilot symbol from pilot sequence $\{\lambda_1, \lambda_2, \dots, \lambda_\tau, \dots, \lambda_{\tau_p}\} \in \mathbb{C}^{\tau \times 1}$ and transmit to the AP cluster for channel estimation, where $\|\lambda_\tau\|^2 = 1$. the pilot signals are orthogonal and have a period of t_p . We can separate users into τ_p sets $S_1, S_2, \dots, S_{\tau_p}$, Users in set S_τ can transmit λ_τ pilot symbol. Let $I_m, \forall m \in \mathcal{M}$ represents the index of pilot symbol transmitted by user m , and we have $I_m \in \{1, 2, \dots, \tau_p\}$. We assume that the receiver at the AP uses large scale fading decoding (LSFD) [5] for signal detection. Each AP first estimates the received signal locally, and then transmits the result to the CPU through fronthaul. The CPU linearly integrates all received local estimates through LSFD vectors to recover the data. Therefore, the signal received by AP n can be given by

$$y_n = \sum_{m=1}^N \sqrt{p_m^d} g_{m,n} x_m + z_m, \quad (6)$$

where x_m is the signal transmit by user m , p_m^d is the payload symbol transmit power of user m , and we have $p_{\min} \leq p_m^d \leq p_{\max}$. $z_m \sim \mathcal{CN}(0, 1)$ is the AWGN. According to [1], AP n can estimate the channel by $\hat{g}_{m,n} \sim \mathcal{CN}(0, \gamma_{m,n})$, where $\gamma_{m,n} = \frac{t_p p_m^p \beta_{m,n}^2}{1 + t_p p_m^p \beta_{m,n}}$, and p_m^p is the power of user m when transmitting pilot symbol. m (We assume that the user sends the pilot signal with the maximum transmit power, i.e., $p_m^p = p_{\max}$)

In UCMEC, each user will be served by a specific AP cluster. In order to save the backhaul bandwidth, we assume that the cluster size is s , and the AP cluster is divided according to the channel quality $g_{m,n}$ in descending order (i.e., take the first s APs with good channel quality). For user m , its AP cluster is represented by Ω_m . For AP n , the set of all users is represented by Φ_n .

The user's access to the AP clusters are coordinated by adopting the orthogonal frequency division multiple access (OFDMA) technique[6]. After each AP performs local estimation, the CPU uses the LSFD vector $v_m = [v_{m,1}, v_{m,2}, \dots, v_{m,N}]^T$ for decoding, thus the SINR of user m is

$$SINR_m = \frac{p_m^d \left(\sum_{n \in \Omega_m} v_{m,n} \gamma_{m,n} \right)^2}{\sum_{i=1}^M p_m^d \sum_{n \in \Omega_m} v_{m,n}^2 \gamma_{m,n} \beta_{i,n} + \sum_{n \in \Omega_m} v_{m,n}^2 \gamma_{m,n}}. \quad (7)$$

To ensure that each user has the maximum SINR, the optimal LSFD vector[7] of user m is given by

$$\mathbf{v}_m = \left(p_m^p \sum_{i \in S_{b_m} \setminus \{m\}} \gamma_m \gamma_m^H + \Lambda_m \right)^{-1} \gamma_m, \quad (8)$$

where

$$\Lambda_m = \text{diag} \left(p_m^p \sum_{m=1}^M \gamma_m \beta_m + \gamma_m \right) \quad (9)$$

where $\gamma_m = \{\gamma_{m,1}, \gamma_{m,2}, \dots, \gamma_{m,n}\}, \forall m \in \mathcal{M}$, $\beta_m = \{\beta_{m,1}, \beta_{m,2}, \dots, \beta_{m,n}\}, \forall m \in \mathcal{M}$, $\mathbf{v}_m = \{v_{m,1}, v_{m,2}, \dots, v_{m,n}\}, \forall m \in \mathcal{M}$.

Therefore, the reachable rate of user m is

$$r_m = B_0 \log_2(1 + \text{SINR}_m) \quad (10)$$

where B_0 is the transmission bandwidth.

1.3 Offloading Model

We assume the computing task $Task_m$ of user m , $m \in \mathcal{M}$ can be represented by a triple, i.e., $Task_m = \{D_m, \rho_m, \tau_m\}$, which represent the data size (in bits), the computing density (in CPU cycles per bit) of the task and the maximum tolerable delay (in second) of the task, respectively. Users can choose to calculate the tasks locally or offloading the tasks to the MEC server on the CPU through AP clusters. We consider the partition offloading ion this paper, and use a continuous variable $a_m \in [0, 1]$ to represent the offloading decision of user m . The a_m part of the task will be offloaded to MEC, and the remaining $1 - a_m$ part of the task is processed locally. Since the downlink data is relatively small, we ignore the overhead of downlink transmission process.

1.3.1 Local Computing Model

We can obtain the delay of user m for local processing by

$$T_m^l = (1 - a_m) \frac{D_m \rho_m}{C_m^l} \quad (11)$$

1.3.2 Edge Computing Model

After the user sends an offloading request, the user will first send data signals to the AP cluster. Therefore, the transmission delay of user m is

$$T_m^t = a_m \frac{D_m}{r_m} \quad (12)$$

We do not consider the overhead of fronthaul transmission in this paper. We assume that the computing resource provided by the MEC server for user m is c_m , thus the computing delay of user m when offloading to the MEC is

$$T_m^e = a_m \frac{D_m \rho_m}{c_m} \quad (13)$$

Thus the offloading delay of user m can be given by

$$T_m^o = T_m^t + T_m^e. \quad (14)$$

Thus the total delay of user m is

$$T_m = \max \{T_m^l, T_m^o\}. \quad (15)$$

1.4 Performance Analysis

1.4.1 Energy Consumption Analysis

The energy consumption of user m when computing the task locally is

$$E_m^l = \eta(1 - a_m)D_m\rho_m (C_m^l)^2 \quad (16)$$

where η is the effective switched capacitance

Since the user needs to send signals to its AP cluster, the energy consumption of the user m for data transmission is

$$E_m^t = \sum_{n \in \Omega_m} p_m^d T_m^t \quad (17)$$

The energy consumed by the user m when waiting for the MEC to process its task is

$$E_m^e = p_m^i T_m^e \quad (18)$$

where p_m^i is the operating power of the user in the idle state. Therefore, the total energy consumed by user m is

$$E_m = E_m^l + E_m^t + E_m^e \quad (19)$$

1.4.2 Successful Transmission Probability (STP) Analysis

According to [8], the successful transmission of user m means that the uplink SINR of the user is greater than a threshold ζ . Thus the STP can be defined as the proportion of users with successful transmission.

1.4.3 Successful Computing Probability (SCP) Analysis

According to [8], the successful computation of user m means that the total delay is less than the delay threshold of the calculation task τ_m . Thus the SCP can be defined as the proportion of users with successful computation.

In addition, we can also define the user with successful offloading as the user who satisfies both successful transmission and computation.

1.5 Problem Formulation

Our optimization goal is to minimize the total energy consumption of all users in the network. Assume that the offloading decision matrix of all users is $\mathbf{A} = \{a_m, m \in \mathcal{M}\}$; The transmit power

decision matrix of all users is $\mathbf{P} = \{p_m^d, m \in \mathcal{M}\}$. The computing resource allocation matrix is $\mathbf{C} = \{c_m, m \in \mathcal{M}\}$. The optimization problem can be formulated as

$$\text{P1 : } \min_{\mathbf{A}, \mathbf{P}, \mathbf{C}} \sum_{m=1}^M E_m \quad (20a)$$

$$s.t. \quad \sum_m^{\mathcal{M}} c_m \leq C_e, \quad (20b)$$

$$p_{\min} \leq p_m^d \leq p_{\max}, \forall m \in \mathcal{M}, \quad (20c)$$

$$T_m \leq \tau_m, \forall m \in \mathcal{M}, \quad (20d)$$

$$0 \leq c_m \leq C_e, \quad (20e)$$

$$a_m \in [0, 1], \forall n \in \mathcal{N}, \forall m \in \mathcal{M}, \quad (20f)$$

where (20b) represents that the amount of computing resource allocated by the MEC server cannot exceed the computing resource capacity. (20c) represents the transmit power of the user should be within a range; (20d) represents the total delay of the user cannot exceed the maximum tolerable delay of the task. (20e) indicating the the feasible region of allocated computing resource. (20f) represents the feasible region of the offloading decision variable. Problem (20) is a nonlinear problem with coupling between variables. In addition, the optimization objective is non convex and the problem can be proved to be NP-hard.

2 Algorithm Design

In this paper, we use the proximal policy optimization (PPO)-based DRL algorithm to solve the above optimization problem. First, we give the Markov decision process (MDP) form of the optimization problem, and then give the derivation of PPO-based DRL algorithm and the corresponding pseudo code.

We can divide a coherent time into several discrete time slots $\mathcal{T} = \{1, 2, \dots, t, \dots, T\}$ and the length of each slot is much less than the coherence time t_c . Each time slot is the minimum time unit for agent to make decisions. At the beginning of each time slot, the agent observes the environment of the ucmec system. We can consider the user's offloading process in ucmec as an MDP, which includes the following elements $\{\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P}, \gamma\}$. In time slot t , $\mathcal{S} = \{\mathbf{s}(t), \forall t \in \mathcal{T}\}$ is the state set, $\mathcal{A} = \{\mathbf{a}(t), \forall t \in \mathcal{T}\}$ is the action set; $\mathcal{P} = \{P(s(t+1)|s(t), a(t)), \forall t \in \mathcal{T}\}$ is the state transition probability set, which can describe the dynamic changes of the system. $\mathcal{R} = \{r(t+1)|s(t), a(t), \forall t \in \mathcal{T}\}$ is the reward set, which can be used to describe the reward received by the agent for doing a certain action in a certain state. $\gamma \in [0, 1]$ is the future discount factor, which is used for describing the weight of future reward to the total reward. A larger γ means that the agent is more concerned about the reward in the future, otherwise it can be considered short-sighted. In time slot t , the agent obtains the state set by observing the environment $\mathbf{s}(t)$, and make an action $\mathbf{a}(t)$, then the environment will move to the next state according to the state transition probability $s(t+1)$, and give corresponding rewards to the agent $R(t)$. The above process will be repeated until the reward coverge to a relatively maximum value. At this time,

we can regard the optimal strategy is obtained. Next, we will give the definitions of observation set, action set and reward set, and give the policy optimization objective of the agent.

- **State Space.** The state space in time slot t can be defined as $\mathbf{s}(t) \triangleq \{r_m(t), \forall t \in \mathcal{T}\}$
- **Action Space.** The action space in time slot t can be defined as $\mathbf{a}(t) \triangleq [a_m(t), p_m^d(t), c_m(t) \forall m \in \mathcal{M}, \forall t \in \mathcal{T}]$
- **Reward Space.** After the agent makes actions, the system will move to the next state $\mathbf{s}(t+1)$. In each time slot, the agent will obtain a reward according to the current state and the action it taken. We assume the reward in time slot t as $r(t) = -\kappa_1 * \sum_{m=1}^M E_m + \kappa_2 * \sum_{m=1}^M (\tau_m - T_m) + \kappa_3 * C_e - \sum_{m=1}^M (c_m)$, where $\kappa_1, \kappa_2, \kappa_3$ are adjustable parameters.

We assume that the policy of the agent is a variable about the parameter $\boldsymbol{\theta}$, which is represented by $\pi(\boldsymbol{\theta})$. We can get the optimization objectives of agent as follows

$$\begin{aligned} \boldsymbol{\theta}^* &= \operatorname{argmax}_{\boldsymbol{\theta}} L(\pi(\boldsymbol{\theta})) \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E} [V_{\pi(\boldsymbol{\theta})}(\mathbf{s})] \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E} [Q_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a})], \end{aligned} \quad (21)$$

where

$$V_{\pi(\boldsymbol{\theta})}(\mathbf{s}) = \mathbb{E} [R(t) | \mathbf{s}(t) = \mathbf{s}], \quad (22)$$

$$Q_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a}) = \mathbb{E} [R(t) | \mathbf{s}(t) = \mathbf{s}, \mathbf{a}(t) = \mathbf{a}], \quad (23)$$

$$R(t) = \sum_{j=t}^T \gamma^{j-t} r(t). \quad (24)$$

In the above formula, $V_{\pi(\boldsymbol{\theta})}$ is the value function for observation, $Q_{\pi(\boldsymbol{\theta})}$ is the value function for observation and action. $R(t)$ is the discounted expected future award of user m in time slot t .

We assign an actor and a critical neural network to the agent and use PPO algorithm to train the agent. According to [9], the policy gradient of the agent can be calculated by

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} L(\pi(\boldsymbol{\theta})) &= \mathbb{E}_{\pi(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}} \log_{\pi(\boldsymbol{\theta})} Q_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a})] \\ &= \mathbb{E}_{\pi(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}} \log_{\pi(\boldsymbol{\theta})} A_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a})] \\ &\approx \mathbb{E}_{\pi(\hat{\boldsymbol{\theta}})} [f \nabla_{\boldsymbol{\theta}} \log_{\pi(\boldsymbol{\theta})} A_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a})], \end{aligned} \quad (25)$$

where $f = \frac{\pi(\boldsymbol{\theta})}{\pi(\hat{\boldsymbol{\theta}})}$, $A_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a}) = Q_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a}) - V_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a})$ is the advantage function for observation and action, $\hat{\boldsymbol{\theta}}$ is the policy parameter obtained by sampling. According to PPO[10] theory, we can clip the policy gradient as:

$$\nabla_{\boldsymbol{\theta}} L(\pi(\boldsymbol{\theta})) \approx \mathbb{E}_{\pi(\hat{\boldsymbol{\theta}})} [\nabla_{\boldsymbol{\theta}} \log_{\pi(\boldsymbol{\theta})} C(\mathbf{s}, \mathbf{a})], \quad (26)$$

where

$$C(\mathbf{s}, \mathbf{a}) = \min [f A_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a}), v(f) A_{\pi(\boldsymbol{\theta})}(\mathbf{s}, \mathbf{a})], \quad (27)$$

$$v(x) = \begin{cases} 1 + \epsilon, & x > 1 + \epsilon, \\ x, & 1 - \epsilon \leq x \leq 1 + \epsilon, \\ 1 - \epsilon, & x < 1 - \epsilon. \end{cases} \quad (28)$$

where ϵ is an adjustable parameter. We define the loss function of the critic network as

$$J(\phi) = \mathbb{E}_{\pi(\phi)} \left[-V_{\phi}(\mathbf{s}) + \mathbb{E}_{\pi(\phi)} \left[r + V_{\phi}(\mathbf{s}') \right] \right] \quad (29)$$

where ϕ is the parameter of the critic network, \mathbf{s}' is the observation result of the agent at the next time slot. The estimated gradient of the critic network is calculated by

$$\nabla_{\phi} J' = \frac{1}{B} \sum_{t=0}^{B-1} \left[V_{\phi}(\mathbf{s}) - Y(t) \frac{dV_{\phi}(\mathbf{s})}{d\phi} \right], \quad (30)$$

where

$$Y(t) = R(t) - \gamma^{B-t} R(B) + \gamma^{B-t} V_{\phi}(B), \quad (31)$$

B is the mini-batch size when updating the critic network. Therefore, the parameter of the critic network can be updated through mini-batch random gradient descent method

$$\phi \leftarrow \phi - l^c \nabla_{\phi} J' \quad (32)$$

where l^c is the learning rate of updating the critic network. Similarly, the estimated gradient of the actor network is calculated by

$$\nabla_{\theta} L' = \frac{1}{B} \sum_{t=0}^{B-1} \nabla_{\theta} \log_{\pi(\theta)} C(\mathbf{s}(t), \mathbf{a}(t)), \quad (33)$$

Thus θ can be updated through mini-batch random gradient descent method

$$\theta \leftarrow \theta + l^a \nabla_{\theta} L', \quad (34)$$

where l^a is the learning rate of updating the actor network.

In summary, the PPO-based offloading and resource allocation algorithm is summarized in **Algorithm 1**.

Algorithm 1 PPO-based Offloading and Resource Allocation Algorithm

Initialization:Agent initialization $\gamma, l^a, l^c, \theta, \phi, \forall m \in \mathcal{M}$ **for** time slot t in $1, 2, \dots, T$ **do**Observe from the environment, update $\mathbf{s}(t)$.Input $\mathbf{s}(t)$ to the actor network, output the actions in current time slot $\mathbf{a}(t)$.Calculate current reward $r(t)$.Store $\{\mathbf{s}(t), \mathbf{a}(t), r(t)\}$ to \mathcal{B} .**if** $t \% B == 0$ **then**Calculate $\nabla_{\theta} L'$ and $\nabla_{\theta} J'$.Update θ and ϕ .Clear the replay Buffer \mathcal{B} **end if****end for****Output:** Optimal offloading strategy, power control strategy and computing resource allocation strategy $\mathbf{A}^*, \mathbf{P}^*, \mathbf{C}^*$

3 Simulation Settings

Table 1: SUPER PARAMETERS of PPO

Parameter	Value
Learning rate of actor network l^a	3×10^{-4}
Learning rate of critic network l^c	3×10^{-4}
Discount factor γ	0.99
GAE parameter	0.95
PPO clip parameter	0.2
Batch size	2048
Mini batch size	64
Number of hidden fully connected layers (Orthogonal Initialization)	2
Number of neurons in each hidden layer	64

Table 2: SYSTEM SIMULATION PARAMETERS

Parameter	Value
Maximum Transmit power of users p_{max}	100mW
Minimum Transmit power of users p_{min}	0 W
Idle power of users p_{idle}	10mW
Computation resource of MEC server C_e	40 GHz
Computation resource of user m	[0.6,1] GHz
Task data size $D_{L_{n,k}}$	[0.05,0.1] MB
Computation density ρ_n	[10000,18000] CPU cycles/bit
Max tolerance delay τ_m	[1,2]s
effective switched capacitance of users η	10^{-27}
Gaussian white noise power z	-174 dBm/Hz
carrier frequency f	1.9GHz
system bandwidth W	20MHz
standard deviation of the shadow fading σ_s	8dB
d_0	10m
d_1	50m
antenna height at the AP h_a	15m
antenna height at the user h_u	1.65m
δ	0.2
coherence time τ_c	1ms

References

- [1] H. Ngo, A. Ashikhmin, H. Yang, E. Larsson, and T. Marzetta, “Cell-free massive mimo versus small cells,” *IEEE Transactions on Wireless Communications*, vol. PP, Feb. 2016.
- [2] S. Buzzi, C. D’Andrea, and A. Zappone, “User-centric 5g cellular networks: Resource allocation and comparison with the cell-free massive mimo approach,” *IEEE Transactions on Wireless Communications*, vol. PP, Mar. 2018.
- [3] Z. Wang, E. Tameh, and A. Nix, “Joint shadowing process in urban peer-to-peer radio channels,” *IEEE Transactions on Vehicular Technology*, vol. 57, pp. 52 – 64, Feb. 2008.
- [4] Z. Chen, E. Björnson, and E. Larsson, “Dynamic resource allocation in co-located and cell-free massive mimo,” *IEEE Transactions on Green Communications and Networking*, vol. PP, pp. 1–1, Nov. 2019.
- [5] E. Björnson and L. Sanguinetti, “Making cell-free massive mimo competitive with mmse processing and centralized implementation,” *IEEE Transactions on Wireless Communications*, vol. PP, pp. 1–1, Sep. 2019.

- [6] S. Seng, C. Luo, X. Li, H. Zhang, and H. Ji, “User matching on blockchain for computation offloading in ultra-dense wireless networks,” *IEEE Transactions on Network Science and Engineering*, vol. 8, no. 2, pp. 1167–1177, Jun. 2021.
- [7] E. Nayebi, A. Ashikhmin, T. Marzetta, and B. Rao, “Performance of cell-free massive mimo systems with mmse and lsfd receivers,” Nov. 2016, pp. 203–207.
- [8] S. Mukherjee and J. Lee, “Edge computing-enabled cell-free massive mimo systems,” *IEEE Transactions on Wireless Communications*, vol. PP, pp. 1–1, Jan. 2020.
- [9] S. Li, X. Hu, and Y. Du, “Deep reinforcement learning and game theory for computation offloading in dynamic edge computing markets,” *IEEE Access*, vol. PP, pp. 1–1, Aug. 2021.
- [10] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” Jul. 2017.