# Joint Computing, Pushing, and Caching Optimization for Mobile Edge Computing Networks via Soft Actor-Critic Learning

Xiangyu Gao, *Student Member, IEEE*, Yaping Sun, Hao Chen, Xiaodong Xu, Shuguang Cui, *Fellow, IEEE*

*Abstract*—Mobile edge computing (MEC) networks bring computing and storage capabilities closer to edge devices, which reduces latency and improves network performance. However, to further reduce transmission and computation costs while satisfying user-perceived quality of experience, a joint optimization in computing, pushing, and caching is needed. In this paper, we formulate the joint-design problem in MEC networks as an infinite-horizon discounted-cost Markov decision process and solve it using a deep reinforcement learning (DRL)-based framework that enables the dynamic orchestration of computing, pushing, and caching. Through the deep networks embedded in the DRL structure, our framework can implicitly predict user future requests and push or cache the appropriate content to effectively enhance system performance. One issue we encountered when considering three functions collectively is the curse of dimensionality for the action space. To address it, we relaxed the discrete action space into a continuous space and then adopted soft actor-critic learning to solve the optimization problem, followed by utilizing a vector quantization method to obtain the desired discrete action. Additionally, an action correction method was proposed to compress the action space further and accelerate the convergence. Our simulations under the setting of a general single-user, single-server MEC network with dynamic transmission link quality demonstrate that the proposed framework effectively decreases transmission bandwidth and computing cost by proactively pushing data on future demand to users and jointly optimizing the three functions. We also conduct extensive parameter tuning analysis, which shows that our approach outperforms the baselines under various parameter settings.

*Index Terms*—computing, pushing, caching, mobile edge computing network, deep reinforcement learning, soft actor-critic.

## I. INTRODUCTION

Recent advancements in smart mobile devices have enabled various emerging applications such as virtual / augmented real-

X. Gao is with the Department of Electrical and Computer Engineering, University of Washington, Seattle, WA, USA. (email: xygao@uw.edu)

Y. Sun and H. Chen are with the Department of Broadband Communication, Peng Cheng Laboratory, Shenzhen 518000, China. (email: {sunyp, chenh03}@pcl.ac.cn)

X. Xu is with the Beijing University of Posts and Telecommunications, Beijing 100876, China, and affiliated with the Department of Broadband Communication, Peng Cheng Laboratory, Shenzhen 518000, China. (email: xuxiaodong@bupt.edu.cn)

S. Cui is with the School of Science and Engineering (SSE) and the Future Network of Intelligent Institute (FNii), the Chinese University of Hong Kong (Shenzhen), Shenzhen 518172, China. S. Cui is also with Shenzhen Research Institute of Big Data, Shenzhen 518172, China, and affiliated with the Department of Broadband Communication, Peng Cheng Laboratory, Shenzhen 518000, China (email: shuguangcui@cuhk.edu.cn).

Corresponding Author: Yaping Sun.

ity (VR/AR) [1], online gaming, and autonomous driving [2]–[4]. These applications require ultra-high communication and computation requirements, making it challenging for mobile operators to minimize communication and computation costs while ensuring the user-perceived quality of experience [5]. In response to these challenges, the mobile edge computing network has emerged as a promising solution by pushing caching and computing resources to access points, base stations (BSs), and mobile devices at the wireless network edge [1].

### A. Prior Art: Caching, Pushing, and Computing Design

Caching can significantly improve bandwidth utilization by placing popular content closer to users for future reuse, leveraging the high degree of asynchronous content reuse in mobile traffic [6]. Caching policies can be classified into two types, *static caching* and *dynamic caching*, depending on whether the cached contents remain unchanged or are dynamically updated. Static caching policies are generally determined based on the content popularity distribution, and the cache state remains unchanged for a relatively long time [6], [7]. In [7], a collaborative content caching scheme among base stations (BSs) in cache-enabled multi-cell cooperative networks is considered to minimize the average request delay, formulated by the stochastic request traffic modeling. Dynamic caching policies, on the other hand, involve updating the content placement from time to time by using instantaneous user request information, such as the least recently used (LRU) and least frequently used (LFU) policy [8]. However, since most caching policies are reactive operations and do not consider proactive pushing, the system performance can be further improved.

Joint pushing and caching can indeed improve system performance by proactively transmitting contents during low traffic periods to satisfy future user demands. Several studies have explored this approach, such as [9] which considers a multi-user wireless network with multicast opportunities to minimize current and future transmission costs, and [10] which uses content request delay information to predict a user's request time for certain content items and maximize effective throughput. Additionally, [11] builds on RDI to characterize asynchronous user requests and proposes a coded joint pushing and caching method to minimize network traffic load with low complexity. In [12], a continuous-time optimization problem is formulated to determine optimal transmission and caching policies for small cell and Device-to-Device networks with

known user demands in advance. However, existing joint pushing and caching policies only consider content delivery and have not taken into account the computation part, which limits their applicability to modern mobile traffic services such as mobile VR delivery.

To effectively serve modern mobile traffic, the joint design of caching, computing, and communication (3C) has been receiving increasing attention. One direction of 3C research focuses on the joint utilization of cache and computing resources at MEC servers to minimize transmission latency [13], [14] and energy consumption [15]. Another direction of 3C research considers the joint utilization of 3C resources at mobile devices to minimize communication costs in both single-user scenarios [1], [16] and multiple-user scenarios [17]. However, these joint 3C designs consider only static caching, where the cache states remain unchanged and do not take into account the benefits of pushing. Therefore, the system performance can be further improved through dynamic caching policies.

### B. DRL-based Systems

Recent advances in deep learning (DL) have enabled the development of novel approaches for complicated classification and detection tasks [18], [19], as well as the solving of complex optimization problems that traditional methods may not be effective or efficient at handling [20]–[27]. Among all popular DL models, reinforcement learning (RL) [28]–[30] has been widely used in scheduling and optimization problems, such as transportation and resource allocation [31]–[33], by learning an optimal policy for the agent to take actions that maximize a reward signal. By using RL, an agent can learn from experience and adapt its behavior over time to achieve the best possible outcomes. For example, in [20], a hierarchical RL algorithm is proposed to solve the joint optimization of pushing and caching in a multi-access edge computing network with multiuser and multicast data. The objective is to maximize bandwidth utilization and decrease the total quantity of data transmitted. In [21], the actor-critic RL framework is utilized to solve the joint optimization of caching, computation offloading, and radio resource allocation in the fog-enabled Internet of Things (IoTs), with the aim of minimizing the average end-to-end delay. In [22], Ning *et al.* develop an intent-based traffic control system that utilizes DRL for the 5G-envisioned Internet of Connected Vehicles, which can dynamically orchestrate edge computing and content caching to improve the profits of mobile network operators. Furthermore, in [23], a distributed DL-based offloading algorithm is proposed, which uses multiple parallel deep neural networks to generate offloading decisions for MEC networks, where multiple wireless devices choose to offload their computation tasks to an edge server. In [24], [25], Zhao *et al.* and Huang *et al.* devise MEC networks for IoTs by using DRL frameworks to make the offloading strategy for offloading some computational tasks from IoT users to the computational access points or MEC server to reduce system latency and energy consumption. However, a common issue encountered when applying DRL-based systems to real-world optimization problems is the curse of dimensionality, which cannot be effectively and efficiently solved by general frameworks and optimization tools, especially for large-scale networks and tasks.

### C. Contributions

To address the issues mentioned earlier, we propose a joint computing, pushing, and caching policy optimization framework in MEC networks[1]. Our contributions are as follows:

- We propose a model for the MEC network that computes its transmission and computation costs while taking into account computing, pushing, and caching actions. By representing system requests and their transition probabilities through a first-order F-state Markov chain, we formulate the joint optimization problem as an infinite-horizon discounted-cost Markov decision process with the dual objectives of reducing transmission and computation costs. Solving this problem requires dynamically optimizing the computing, pushing, and caching decisions over time to achieve the best overall performance.

- To address the curse of dimensionality in the joint optimization problem, we implemented a continuous-space DRL approach known as soft actor-critic (SAC) learning [29]. Unlike classic discrete-space DRL algorithms, such as deep Q-learning [30], which rely on the Q-networks with a size linearly increased with the action space, SAC only requires learning the Gaussian-format Q-functions [29]. As a result, SAC significantly reduces the number of parameters that need to be learned in a neural network. However, this does introduce the challenge of having an output action in continuous space that cannot be directly utilized. Therefore, we have designed an action quantization and correction algorithm that allows us to tailor SAC to our discrete optimization problem. Furthermore, the SAC algorithm is known for its stability and ease of convergence [29].

- We present simulation results with various system parameters under the setting of a general single-user, single-server MEC network to demonstrate the effectiveness of the proposed SAC algorithm. Our results show that by considering the joint optimization of computing, pushing, and caching, the performance of the MEC network can be significantly improved in terms of lower computation cost and reduced transmission cost. Moreover, our approach outperforms baseline methods that consider only a subset of these functions, demonstrating the benefits of the joint optimization.

### D. Outline

The paper is organized as follows: Section II outlines the system model for the MEC network. Section III formulates the joint policy optimization problem. Section IV presents the utilization of SAC in optimization. Section V covers implementation details and evaluation results. Section VI concludes the paper.

---

[1]The code and sample data of this framework will be made open-source and available at *https://github.com/Xiangyu-Gao/sac_joint_compute_push_cache*
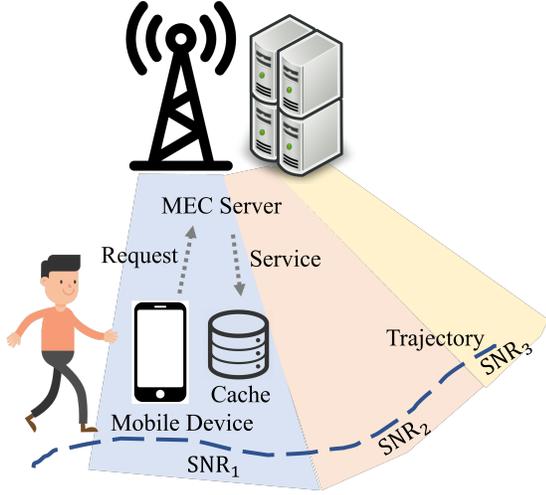
Fig. 1: Illustration of MEC network with single MEC server and single mobile device. The mobile device is assumed to be moving at a small speed, such as an iPhone being carried by an individual. The channel quality for communication between the mobile device and the MEC server is modeled as an SNR, which may change over time due to the movement of the mobile device.

## II. SYSTEM MODEL

Without loss of generality, we begin by considering a simple mobile edge network consisting of one MEC server and one mobile device, as shown in Fig. 1. The system model can be extended to the multi-user scenario by summing the objective functions of multiple users and considering the restrictions of the total communication and computing resources. The MEC server has a large cache size, sufficient to proactively store the input and output data of all tasks requested by the mobile device. In contrast, the cache size of the mobile device is limited to a capacity denoted as $C$ (in bits). The mobile device is equipped with multi-core computing capabilities, each with a computation frequency $f_D$ (in cycles/s), and the number of computing cores is assumed to be $M$. The system operates over an infinite time horizon, with time slotted and indexed by $t = 0, 1, 2, \cdots$, each with a fixed length of $\tau$ seconds. At the start of each time slot, the mobile device submits one task request, which is delay-intolerant and must be served before the end of the slot. The tasks are categorized as delay-intolerant due to the critical importance of upholding optimal user experience and ensuring high-quality service for applications such as AR/VR, real-time communication, and streaming applications which are notably sensitive to delays. Due to the mobility of the device, the data transmission rate for the link between the mobile device and the server may vary over time. To model this dynamic effect, we adopt the signal-to-noise ratio (SNR) following the Shannon theory, which measures the quality of the link. The system is designed to optimize the joint pushing, caching, and computing functions to minimize the computation and transmission costs of the network while ensuring timely and efficient task execution.

### A. Task Model

Assuming that the mobile device will request a total of $F$ tasks, we define the task set $\mathcal{F}$ as $\mathcal{F} \triangleq \{1, 2, \ldots, f, \ldots, F\}$. Each task $f \in \mathcal{F}$ is characterized by a 4-item tuple $\left\{ I_f \text{ (in bits)}, O_f \text{ (in bits)}, w_f \text{ (in cycles/bit)}, \tau \text{ (in seconds)} \right\}$. Specifically, $I_f$ represents the size of the input data generated from the Internet which can be cached. $O_f$ represents the size of the output data after the computation is completed[2]. $w_f$ and $\tau$ denote the required computation cycles per bit and the maximum service latency, respectively.

### B. System State

*1) Request State:* At each time slot $t$, the mobile device submits a single task request. The request state at time $t$ is denoted by $A(t) \in \mathcal{F}$ representing the requested task, where $A(t) = f$ signifies that task $f$ in set $\mathcal{F}$ is being requested by the mobile device. The size of $\mathcal{F}$ is $F$. To model the evolution of requested tasks and their transition probabilities, we employed a first-order F-state Markov chain [20], [34], referred to as $A(t) : t = 0, 1, 2, \cdots$. In this context, each state within the Markov chain corresponds to a distinct task, and the total number of tasks is assumed to be $F$. The choice to use a first-order Markov chain is rooted in its assumption that the probability of transitioning to a particular state is solely dependent on the current state. The probability of transitioning to state $j \in \mathcal{F}$ at time slot $t+1$, given that the request state at time slot $t$ is $i \in \mathcal{F}$, is represented by $\Pr[A(t+1) = j | A(t) = i]$. It is assumed that $A(t)$ is time-homogeneous. We denote the transition probability matrix of $A(t)$ with $\mathbf{Q} \triangleq (q_{i,j})_{i \in \mathcal{F}, j \in \mathcal{F}}$, where $q_{i,j} \triangleq \Pr[A(t+1) = j | A(t) = i]$. Moreover, we focus our attention on an irreducible Markov chain to reflect the idea that any state in the system can be reached from any other state with a non-zero probability. We denote the limiting distribution of $A(t)$ with $\mathbf{p} \triangleq (p_f)_{f \in \mathcal{F}}$. Here, $p_f \triangleq \lim_{t \to \infty} \Pr[A(t) = f]$, and it should be noted that $p_f = \sum_{i \in \mathcal{F}} p_i q_{i,f}$ for all $f \in \mathcal{F}$.

*2) Cache State:* Let $S_f^I(t) \in \{0, 1\}$ denote the indicator of the cache state of the input data for task $f$ stored in the mobile device. Here, $S_f^I(t) = 1$ means that the input data for task $f$ is cached in the mobile device, while $S_f^I(t) = 0$ implies that the input data is not cached. Similarly, let $S_f^O(t) \in \{0, 1\}$ denote the indicator of the cache state of the output data for task $f$ stored in the mobile device, where $S_f^O(t) = 1$ represents that the output data for task $f$ is cached in the mobile device, and $S_f^O(t) = 0$ implies that the output data is not cached. The cache size of the mobile device is denoted by $C$ (in bits). The cache size constraint is given by

$$\sum_{f=1}^{F} I_f S_f^I(t) + O_f S_f^O(t) \leq C \tag{1}$$

which enforces that the sum of the sizes of input and output data cached for all tasks in the mobile device cannot exceed the cache size.

---

[2]In many systems, the actual output size might not be known beforehand, especially for computational tasks that involve dynamic data processing. It's possible that we could use historical data or estimations based on the characteristics of the input data and the computation process.

We define the cache state of the mobile device at time slot $t$, denoted by $\mathbf{S}(t) \triangleq (S_f^I(t), S_f^O(t))_{f \in \mathcal{F}} \in \mathcal{S}$, where $\mathcal{S} \triangleq \{(S_f^I, S_f^O)_{f \in \mathcal{F}} \in \{0,1\}^F \times \{0,1\}^F : \sum_{f \in \mathcal{F}} I_f S_f^I + O_f S_f^O \leq C\}$ represents the cache state space of the mobile device. Here, $N_{\min} \triangleq \frac{C}{\max_{f \in \mathcal{F}} I_f, O_f}$ and $N_{\max} \triangleq \frac{C}{\min_{f \in \mathcal{F}} I_f, O_f}$ represent the lower and upper bounds, respectively, on the cardinality of $\mathcal{S}$. The cardinality of $\mathcal{S}$ is bounded by $\binom{F}{N_{\min}}$ and $\binom{F}{N_{\max}}$ from below and above, respectively.

*3) System State:* At time slot $t$, the system state consists of both system request state and system cache state, represented by $\mathbf{X}(t) \triangleq (A(t), \mathbf{S}(t)) \in \mathcal{F} \times \mathcal{S}$, where $\mathcal{F} \times \mathcal{S}$ represents the system state space.

## C. System Action

*1) Reactive Computation Action:* At each time slot $t$, the reactive transmission bandwidth cost and the reactive computation energy cost are denoted as $B^R(t)$ and $E^R(t)$, respectively. The task request $A(t)$ is served based on the current system state $\mathbf{X}(t) = (A(t), \mathbf{S}(t))$ as follows:

- If the cache state $S_{A(t)}^O(t)$ is equal to 1, it indicates that the output of task $A(t)$ is already cached locally, hence it can be retrieved without the need for any transmission or computation. As a result, the delay is negligible, and both the reactive computation energy and transmission cost become zero.
- Assuming that $S_{A(t)}^I(t) = 1$ and $S_{A(t)}^O(t) = 0$, it is possible to compute the requested task $A(t)$ directly using the locally cached input data. Let us define $c_{R,f}(t) \in \{1, \cdots, M\}$ as the number of computation cores allocated for reactively processing task $f$ at time slot $t$ on the mobile device. Consequently, we can set $c_{R,f}(t) = 0$ for all $f \in \mathcal{F} \backslash A(t)$. To ensure that the requested task $A(t)$ is completed within $\tau$, we must have $\frac{I_{A(t)} w_{A(t)}}{\tau} \leq c_{R,A(t)}(t) f_D$.[3] Here $I_{A(t)}$ and $w_{A(t)}$ denote the input size and the computational workload of task $A(t)$, respectively. We can calculate the energy consumed for computing one cycle with frequency $c_{R,f}(t) f_D$ on the mobile device as $\mu c_{R,f}^2(t) f_D^2$, where $\mu$ is the effective switched capacitance related to the chip architecture indicating the power efficiency of the CPU. Therefore, the reactive computation energy cost $E^R(t)$ is given by $\mu c_{R,A(t)}^2(t) f_D^2 I_{A(t)} w_{A(t)}$, and the reactive transmission cost $B^R(t)$ is zero.
- If $S_{A(t)}^I(t) = 0$ and $S_{A(t)}^O(t) = 0$, the mobile device must download the input data of task $A(t)$ from the MEC server before computing it locally. Let $\text{SNR}(t)$ be the SNR value of the data transmission link at time slot $t$. The required latency can be expressed as $\frac{I_{A(t)}}{B^R(t) \log_2(1+\text{SNR}(t))} + \frac{I_{A(t)} w_{A(t)}}{c_{R,A(t)}(t) f_D}$, where $B^R(t) \log_2(1 + \text{SNR}(t))$ is the channel capacity given by Shannon theory. To satisfy the latency constraint, i.e., $\frac{I_{A(t)}}{B^R(t) \log_2(1+\text{SNR}(t))} + \frac{I_{A(t)} w_{A(t)}}{c_{R,A(t)}(t) f_D} \leq$

$\tau$, the minimum reactive transmission cost $B^R(t)$ is given by $\frac{I_{A(t)}}{\left(\tau - \frac{I_{A(t)} w_{A(t)}}{c_{R,A(t)}(t) f_D}\right) \log_2(1+\text{SNR}(t))}$.[4] The reactive computation energy cost $E^R(t)$ is given by $\mu c_{R,A(t)}^2(t) f_D^2 I_{A(t)} w_{A(t)}$.

In summary, at time slot $t$, the reactive computation action $c_{R,f}(t)$ should satisfy

$$c_{R,f}(t) \leq \mathbf{1}(A(t) = f)\left(1 - S_f^O(t)\right) M, \ \forall f \in \mathcal{F}, \quad (2)$$

and the reactive transmission cost $B^R(t)$ is given by

$$B^R(t) = \left(1 - S_{A(t)}^I(t)\right)\left(1 - S_{A(t)}^O(t)\right) \quad (3)$$

$$\times \frac{I_{A(t)}}{\left(\tau - \frac{I_{A(t)} w_{A(t)}}{c_{R,A(t)}(t) f_D}\right) \log_2(1 + \text{SNR}(t))}, \quad (4)$$

and the reactive computation cost $E^R(t)$ is given by

$$E^R(t) = \left(1 - S_{A(t)}^O(t)\right) \mu c_{R,A(t)}^2(t) f_D^2 I_{A(t)} w_{A(t)}. \quad (5)$$

Let $\mathbf{c}_R \triangleq (c_{R,f})_{f \in \mathcal{F}} \in \Pi_C^R(\mathbf{X})$ denote the reactive computation action of the system, where $\Pi_C^R(\mathbf{X}) \triangleq \left\{(c_{R,f})_{f \in \mathcal{F}} \in \{0, 1, \cdots, M\}^F : (2)\right\}$ represents the decision space for reactive computation of the system under state $\mathbf{X}$. It can be observed from Eq. (2) that the size of the reactive computation action space is $M + 1$.

*2) Proactive Transmission or Pushing Action:* Let $b_f(t) \in \{0, 1\}$ denote the binary decision variable for task $f \in \mathcal{F}$, where $b_f(t) = 1$ indicates that the remote input data of task $f$ is pushed to the mobile device, and $b_f(t) = 0$ otherwise. We assume that the pushed data is transmitted to the mobile device by the end of the time slot. To ensure compliance with the latency constraint, we enforce $\frac{\sum_{f=1}^F I_f b_f(t)}{\tau} \leq B^P(t) \log_2(1 + \text{SNR}(t))$, where $B^P(t)$ denotes the proactive transmission bandwidth cost. Thus, the minimum proactive transmission cost can be expressed as:

$$B^P(t) = \frac{\sum_{f=1}^F I_f b_f(t)}{\tau \log_2(1 + \text{SNR}(t))}. \quad (6)$$

In summary, the system pushing action under system state $\mathbf{b} \triangleq (b_f)_{f \in \mathcal{F}} \in \{0,1\}^F$. The size of the system pushing action space under system state $\mathbf{X}$ is $2^F$

*3) Cache Update Action:* The cache state of each task $f \in \mathcal{F}$ is updated according to

$$S_f^I(t+1) = S_f^I(t) + \Delta s_f^I(t), \quad (7)$$
$$S_f^O(t+1) = S_f^O(t) + \Delta s_f^O(t), \quad (8)$$

---

[3]We assume that $\frac{I_f w_f}{\tau} \mathbf{1}(A(t) = f) \leq M f_D$, for feasibility, where $\mathbf{1}(A(t) = f)$ is the indicator function that is equal to 1 if $A(t) = f$, and 0 otherwise, and $M$ is the maximum number of computation cores. This assumption holds for all $f \in \mathcal{F}$.

[4]The steps of deriving $B^R(t)$ from the preceding latency constraint are as follows: First, we have $I_{A(t)} / \left(B^R(t) \log_2(1 + \text{SNR}(t))\right) \leq \tau - I_{A(t)} w_{A(t)} / \left(c_{R,A(t)} f_D\right)$. Then, we can get $I_{A(t)} / B^R(t) \leq \left(\tau - I_{A(t)} w_{A(t)} / \left(c_{R,A(t)} f_D\right)\right) \log_2(1 + \text{SNR}(t))$. Finally, we can get $B^R(t) \geq I_{A(t)} / \left(\left(\tau - I_{A(t)} w_{A(t)} / \left(c_{R,A(t)} f_D\right)\right) \log_2(1 + \text{SNR}(t))\right)$.

where $\Delta s_f^I(t) \in \{-1, 0, 1\}$ and $\Delta s_f^O(t) \in \{-1, 0, 1\}$ denote the update action for the cache state of the input and output data of task $f$, respectively. Then, we have $\forall f \in \mathcal{F}$

$$- S_f^I(t) \le \Delta s_f^I(t) \le \min\left\{b_f(t) + c_{R,f}(t), 1 - S_f^I(t)\right\} \quad (9)$$

$$- S_f^O(t) \le \Delta s_f^O(t) \le \min\left\{c_{R,f}(t), 1 - S_f^O(t)\right\}, \quad (10)$$

$$\sum_{f=1}^{F} I_f\left(S_f^I(t) + \Delta s_f^I(t)\right) + O_f\left(S_f^O(t) + \Delta s_f^O(t)\right) \le C, \quad (11)$$

where the left-hand side of Eq. (9) specifies that the removal of the input of task $f$ from the mobile device is only possible if it has been previously cached. On the other hand, the right-hand side of Eq. (9) indicates that the caching of the input of task $f$ into the mobile device is only allowed if it has not been cached before and if it is either proactively transmitted from the MEC server or reactively transmitted, i.e., if $b_f(t) = 1$ or $c_{R,f}(t) > 0$. Similarly, the left-hand side of Eq. (10) states that the output of task $f$ can only be removed from the mobile device if it has been previously cached. On the other hand, the right-hand side of Eq. (10) specifies that the caching of the output of task $f$ into the mobile device is only allowed if it has not been cached before and if it is reactively computed at the mobile device, i.e., if $c_{R,f}(t) > 0$. Finally, Eq. (11) requires that the updated cache state complies with the cache size constraint.

In summary, let $\Delta\mathbf{s} \triangleq \left(\Delta s_f^I, \Delta s_f^O\right) f \in \mathcal{F} \in \Pi_{\Delta s}(\mathbf{X})$ denote the system cache update action, where $\Pi_{\Delta s}(\mathbf{X}) \triangleq \left\{\left(\Delta s_f^I, \Delta s_f^O\right)_{f \in \mathcal{F}} \in \{-1, 0, 1\}^F \times \{-1, 0, 1\}^F : (9), (10), (11)\right\}$. Here, $\Pi_{\Delta s}(\mathbf{X})$ represents the system cache update action space for the given system state $\mathbf{X}$, and it includes tuples of $\Delta s_f^I$ and $\Delta s_f^O$ for each task $f \in \mathcal{F}$, with values in $\{-1, 0, 1\}$ indicating whether to evict, retain, or cache a task's input and output data.

*4) System Action:* The system action at each time slot is a combination of three distinct actions: reactive computation, pushing, and cache update. This combination is represented as $(\mathbf{c}_R, \mathbf{b}, \Delta\mathbf{s}) \in \Pi(\mathbf{X})$, where $\Pi(\mathbf{X})$ is the system action space under the current system state $\mathbf{X}$, $\Pi(\mathbf{X}) \triangleq \Pi_C^R(\mathbf{X}) \times \{0, 1\}^F \times \Pi_{\Delta s}(\mathbf{X})$.

### D. System Cost

At each time slot $t$, the overall system cost is a combination of two components, namely the transmission bandwidth cost and the computation energy cost. The transmission bandwidth cost consists of both proactive and reactive transmission costs and is given by

$$B(t) = B^R(t) + B^P(t), \quad (12)$$

where $B^R(t)$ is given in Eq. (3) and $B^P(t)$ is given in Eq. (6). he reactive computation cost contributes to the computation energy cost only and is given by

$$E(t) = E^R(t), \quad (13)$$

where $E^R(t)$ is given in Eq. (5).

To strike a balance between communication and computation cost, the system cost at time slot $t$ is computed as the weighted sum of transmission bandwidth cost and computation

energy cost, i.e., $B(t) + \lambda E(t)$, where $\lambda$ is a non-negative weighting factor.

## III. PROBLEM FORMULATION

Given an observed system state $\mathbf{X}$, the joint reactive computing, transmission, and caching action, denoted as $(\mathbf{c}_R, \mathbf{b}, \Delta\mathbf{s})$, is determined according to a policy defined as below.

**Definition 1** (Stationary Joint Computing, Pushing and Caching Policy). *A stationary joint computing, pushing, and caching policy $\pi$ is a mapping from system state $\mathbf{X}$ to system action $(\mathbf{c}_R, \mathbf{b}, \Delta\mathbf{s})$, i.e., $(\mathbf{c}_R, \mathbf{b}, \Delta\mathbf{s}) = \pi(\mathbf{X}) \in \Pi(\mathbf{X})$.*

From properties of $\{A(t)\}$ and $\{\mathbf{S}(t)\}$, the induced system state process $\{\mathbf{X}(t)\}$ under policy $\pi$ is a controlled Markov chain. The expected total discounted cost $\phi(\pi)$ is given as:

$$\phi(\pi) \triangleq \limsup_{T \to \infty} \sum_{t=0}^{T-1} \gamma^t \mathbb{E}\left[B(t) + \lambda E(t)\right], \quad (14)$$

where $T$ is the length of the request process, $\gamma$ is the discount factor, $B(t), E(t)$ are the transmission bandwidth cost and computation energy cost at time $t$, and $\lambda$ is the weight balancing two costs.

In this paper, we aim to obtain optimal joint computing, pushing, and caching policy to minimize the sum of infinite horizon discounted system cost, i.e., minimize both the transmission and computation cost, as follows:

**Problem 1** (Joint Computing, Pushing and Caching Policy Optimization).

$$\phi^* \triangleq \min_{\pi} \quad \phi(\pi)$$
$$s.t. \quad \pi(\mathbf{X}) \in \Pi(\mathbf{X}), \quad \forall \mathbf{X} \in \mathcal{F} \times \mathcal{S}.$$

## IV. SOFT ACTOR-CRITIC LEARNING

### A. SAC System State and Action

The system state $\mathbf{x}$ of SAC is designed the match the system state $\mathbf{X}$ in the formulated problem, such that $\mathbf{x} = \mathbf{X} = (A(t), \mathbf{S}(t))$, with a vector size of $2F + 1$.

The SAC algorithm is designed to solve continuous-action problems, whereas the required system action $(\mathbf{c}_R, \mathbf{b}, \Delta\mathbf{s})$ in the formulated problem is discrete. To address this issue, we define the system action of the SAC as the ***continuous version*** of the formulated system action space. This continuous version is denoted as $\mathbf{a} = (\bar{\mathbf{c}}_R, \bar{\mathbf{b}}, \Delta\bar{\mathbf{s}}) \in \bar{\Pi}(\mathbf{X}) \triangleq \bar{\Pi}_C^R(\mathbf{X}) \times [0, 1]^F \times \bar{\Pi}_{\Delta s}(\mathbf{X})$. Here, $\bar{\Pi}_C^R(\mathbf{X}) \triangleq \left\{(c_{R,f})_{f \in \mathcal{F}} \in [0, M]^F : (2)\right\}$, and $\bar{\Pi}_{\Delta s}(\mathbf{X}) \triangleq \left\{\left(\Delta s_f^I, \Delta s_f^O\right)_{f \in \mathcal{F}} \in [-1, 1]^F \times [-1, 1]^F : (9), (10), (11)\right\}$.

As $\bar{\mathbf{c}}_R \triangleq \left\{(\bar{c}_{R,f})_{f \in \mathcal{F}}\right\}$ must always equal zero for $f \in \mathcal{F} \backslash A(t)$, the action space of SAC can be simplified by disregarding the computing cores for non-requested tasks. We can obtain the simplified form of action $\mathbf{a}$ as $\mathbf{a} = (\bar{c}_{A(t)}, \bar{\mathbf{b}}, \Delta\bar{\mathbf{s}})$, with a vector size of $3F + 1$.

## B. SAC Learning

SAC is an off-policy deep reinforcement learning method that maintains the advantages of entropy maximization and stability while offering sample-efficient learning [29]. It operates on an actor-critic framework where the actor is responsible for maximizing expected reward while simultaneously maximizing entropy. The critic evaluates the effectiveness of the policy being followed.

A general form of maximum-entropy RL is given by:

$$J(\pi) = \sum_{t=0}^{T} \mathbb{E}_{(\mathbf{x}_t, \mathbf{a}_t) \sim \rho_\pi} \left[ r(\mathbf{x}_t, \mathbf{a}_t) + \alpha \mathcal{H}(\pi(\cdot \mid \mathbf{x}_t)) \right] \quad (15)$$

where the temperature parameter $\alpha$ determines the relative importance of the entropy term against the reward $r$, and the entropy term is given by $\mathcal{H}(\pi(\cdot \mid \mathbf{x}_t)) = \mathbb{E}_{\mathbf{a}_t}[-\log \pi(\mathbf{a}_t \mid \mathbf{x}_t)]$.

The SAC algorithm is a policy iteration approach designed to solve the optimization problem in Eq. (15) [29]. It comprises two essential components: soft Q-function $Q_\theta(\mathbf{x}_t, \mathbf{a}_t)$, and policy $\pi_\phi(\mathbf{a}_t \mid \mathbf{x}_t)$. To deal with the large continuous domains, neural networks are utilized to approximate these components, with the network parameters denoted by $\theta$ and $\phi$. For example, the policy is modeled as a Gaussian distribution with a fully connected network providing the mean and covariance value, and the Q-function is also approximated using a fully connected neural network. Following [29], the update rules for $\theta$ and $\phi$ are provided below.

The soft Q-function parameters can be trained to minimize the soft Bellman residual

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{x}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[ \frac{1}{2} \left( Q_\theta(\mathbf{x}_t, \mathbf{a}_t) - \left( r(\mathbf{x}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{x}_{t+1} \sim p}[V_{\bar\theta}(\mathbf{x}_{t+1})] \right) \right)^2 \right], \quad (16)$$

where $\mathcal{D}$ is the distribution of previously sampled states and actions, $p$ is the transition probability between states, and the value function $V_{\bar\theta}(\mathbf{x}_t)$ is implicitly parameterized through the soft Q-function parameters as follows

$$V_{\bar\theta}(\mathbf{x}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi} [Q_{\bar\theta}(\mathbf{x}_t, \mathbf{a}_t) - \alpha \log \pi(\mathbf{a}_t \mid \mathbf{x}_t)] \quad (17)$$

The update makes use of a target soft Q-function $Q_{\bar\theta}$ with parameters $\bar\theta$ obtained as an exponentially moving average of the soft Q-function weights $\theta$, which helps stabilize training. The soft Bellman residual $J_Q(\theta)$ in Eq. (16) can be optimized with stochastic gradients

$$\hat{\nabla}_\theta J_Q(\theta) = \nabla_\theta Q_\theta(\mathbf{a}_t, \mathbf{x}_t) \left( Q_\theta(\mathbf{x}_t, \mathbf{a}_t) - \left( r(\mathbf{x}_t, \mathbf{a}_t) + \gamma \left( Q_{\bar\theta}(\mathbf{x}_{t+1}, \mathbf{a}_{t+1}) - \alpha \log(\pi_\phi(\mathbf{a}_{t+1} \mid \mathbf{x}_{t+1})) \right) \right) \right). \quad (18)$$

The policy parameters $\phi$ can be learned by directly minimizing the expected KL divergence in

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{x}_t \sim \mathcal{D}} \left[ \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} \left[ \alpha \log(\pi_\phi(\mathbf{a}_t \mid \mathbf{x}_t)) - Q_\theta(\mathbf{x}_t, \mathbf{a}_t) \right] \right] \quad (19)$$

A neural network transformation is used to parameterize the policy as $\mathbf{a}_t = f_\phi(\epsilon_t; \mathbf{x}_t)$, where $\epsilon_t$ is an input noise vector sampled from a Gaussian distribution. The objective stated by Eq. (19) can be rewritten as:

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{x}_t \sim \mathcal{D}, \epsilon_t \sim \mathcal{N}} \left[ \alpha \log \pi_\phi(f_\phi(\epsilon_t; \mathbf{x}_t) \mid \mathbf{x}_t) - Q_\theta(\mathbf{x}_t, f_\phi(\epsilon_t; \mathbf{x}_t)) \right], \quad (20)$$

where $\pi_\phi$ is defined implicitly in terms of $f_\phi$. The gradient of Eq. (20) is approximated with

$$\hat{\nabla}_\phi J_\pi(\phi) = \nabla_\phi \alpha \log(\pi_\phi(\mathbf{a}_t \mid \mathbf{x}_t)) + (\nabla_{\mathbf{a}_t} \alpha \log(\pi_\phi(\mathbf{a}_t \mid \mathbf{x}_t)) - \nabla_{\mathbf{a}_t} Q(\mathbf{x}_t, \mathbf{a}_t)) \nabla_\phi f_\phi(\epsilon_t; \mathbf{x}_t), \quad (21)$$

where $\mathbf{a}_t$ is evaluated using $f_\phi(\epsilon_t; \mathbf{x}_t)$.

**Remark**: In the maximum entropy framework, the soft policy iteration that alternates between the policy evaluation Eq. (16) and the policy improvement Eq. (19) converges to the optimal policy. *Proof in* [29].

## C. Action Quantization and Correction

In the context of SAC learning, the output at any given time $t$ corresponds to the SAC action $\mathbf{a}_t$, which seeks to maximize the policy value $\pi_\phi(\mathbf{a}_t \mid \mathbf{x}_t)$ with respect to the current SAC state $\mathbf{x}_t$. In order to assess the reward and update the cache, it is necessary to **discretize the continuous SAC action** $\mathbf{a}_t$ and obtain a discrete action $(c_{A(t)}, \mathbf{b}, \Delta\mathbf{s})$. To achieve this goal, a simple action quantization approach was implemented that relies on thresholding and integer projection.

**Action quantization**: Let us consider an element $\bar\eta$ in the SAC action $\mathbf{a}$ and its corresponding quantized version $\eta$ with the selection set $S_\eta$. To obtain $\eta$ from $\bar\eta$, we adopt a uniform thresholding method for integer projection. Specifically, we use the following equation:

$$\eta = \min S_\eta + (\bar\eta - \min S_\eta) \bmod \frac{\max S_\eta - \min S_\eta}{\max S_\eta - \min S_\eta + 1} \quad (22)$$

As an example, consider the push action $b_f(t) \in S_{b_f} = \{0, 1\}$, we can determine its quantized value $b_f(t) = \bar{b}_f(t)$ mod 0.5 using Eq. (22).

**Action correction**: The valid action space of the system is highly constrained due to the limitations imposed by Eq. (2), (3), (9), (10), and (11), resulting in a sparsely-spanning space with a cardinality of $(M + 1) \times 2^F \times 3^{2F}$. Consequently, even with techniques such as penalty reward, it becomes challenging for the SAC algorithm to identify which actions are valid in this vast space. Therefore, the post-quantization action $(c_{A(t)}, \mathbf{b}, \Delta\mathbf{s})$ obtained from SAC is often invalid. In order to address this issue, we propose *Rules 1, 5, and 7* to ensure that the output action of SAC is valid, while satisfying the constraints outlined in Section II-C. Additionally, we introduce *Rules 2, 3, 4, and 6* to improve the training process and enhance the system's overall performance by further compressing the action space, reducing unnecessary costs, and minimizing waste.

- *Rule 1*: When $S_{A(t)}^O$ equals 0, the system checks if the suggested number of computation cores, denoted as $c_{A(t)}$, is less than the minimum workable value given by $\lceil I_{A(t)} w_{A(t)} / (\tau f_D) \rceil$ where $\lceil \cdot \rceil$ represents rounding up to the nearest integer. If this is the case, $c_{A(t)}$ is updated to $\lceil I_{A(t)} w_{A(t)} / (\tau f_D) \rceil$. On the other hand, if $S_{A(t)}^O$ equals

1, $c_{A(t)}$ is set to 0. These rules are designed to fulfill the service latency constraint and reduce unnecessary computation.

- **Rule 2**: When $S_f^I + S_f^O \geq 1$, we set $b_f = 0$. This rule indicates that there is no need for proactive pushing if any data of a task is already cached.
- **Rule 3**: To minimize the cost of pushing data to the mobile device, we ensure that at most one task is proactively transmitted, and this task must have the largest $\bar{b}_f$ value among all un-pushed tasks. The selected task will have a $b_f$ value of 1, while all other tasks will have a $b_f$ value of 0. This approach is adopted to avoid unnecessary pushing costs, as the mobile device is only capable of processing one task request per time slot.
- **Rule 4**: If $b_f = 1$, we set $\Delta s_f^I = 1$, indicating that the data being proactively pushed needs to be cached.
- **Rule 5**: If the sum of cache sizes given by Eq. (11) exceeds the cache capacity, we drop the input or output cache depending on the ascending order of their corresponding $\bar{s}$ values until the cache capacity is satisfied.
- **Rule 6**: If the sum of the caches given by Eq. (11) is less than the capacity, we attempt to add reactive input or output cache based on the decreasing order of the continuous variables $\Delta \bar{s}_{A(t)}^I$ and $\Delta \bar{s}_{A(t)}^O$.
- **Rule 7**: The cache action $\Delta \mathbf{s}$ should be clipped according to the minimum and maximum limits specified in Eq. (9) and Eq. (10).

### D. Reward Design

The reward function $r(\mathbf{x}, \mathbf{a})$ for the SAC state $\mathbf{x}$ and action $\mathbf{a}$ is defined as a function of the resulting bandwidth and computation cost. Specifically, it is given by

$$r(\mathbf{x}, \mathbf{a}) = -\kappa(B(t) + \lambda E(t)) \tag{23}$$

where $\kappa$ is a normalization coefficient that is set to $10^{-6}$ in this paper.

The complete SAC learning algorithm is presented in Algorithm 1. The step sizes for stochastic gradient descent, $\lambda_Q$, and $\lambda_\pi$ are set to $1 \times 10^{-4}$. The target smoothing coefficient, $\xi$, is chosen to be 0.005.

### V. IMPLEMENTATION AND EVALUATION

#### A. Baselines

The proposed system is built on the ***proactive transmission and dynamic-computing-frequency reactive service with cache***, referred to as **PTDFC**. For comparison, we have selected the following baselines:

- ***Most-recently-used proactive transmission and least-recently-used cache replacement*** (**MRU-LRU**): This is a heuristic algorithm [8], [9], where at each time slot, the requested task is reactively served, and the input data of the most-recently-used task is proactively cached. When the cache is full, the input data cache of the least-recently-used task is replaced. We choose to cache only the input data, excluding the output data (post-calculation), due to the common scenario where output data tends to be

---

**Algorithm 1** SAC Learning for Our Problem

---

Initialize parameters $\theta, \bar{\theta}, \phi$ for networks $Q_\theta, Q_{\bar{\theta}}, \pi_\phi$.
Initialize learning rate $\lambda_Q, \lambda_\pi$, and weight $\xi$.
**for** each iteration **do**
    **for** each environment step **do**
        $\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t \mid \mathbf{x}_t)$
        $\mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{a}_t)$
        $\mathbf{a}_t$ quantization & correction, $r(\mathbf{x}_t, \mathbf{a}_t)$ calculation
        $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{x}_t, \mathbf{a}_t, r(\mathbf{x}_t, \mathbf{a}_t), \mathbf{x}_{t+1})\}$
    **end for**
    **for** each gradient step **do**
        $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$ for $i \in \{1, 2\}$
        $\phi \leftarrow \phi - \lambda_\pi \nabla_\phi J_\pi(\phi)$
        $\bar{\theta}_i \leftarrow \xi \theta_i + (1 - \xi)\bar{\theta}_i$ for $i \in \{1, 2\}$
    **end for**
**end for**

---

larger in size than input data. This size difference makes caching output data less efficient in the heuristic design for the purpose of reducing overall costs. The number of computing cores being used is fixed at $0.75M$.

- ***Most-frequently-used proactive transmission and least-frequently-used cache replacement*** (**MFU-LFU**): This algorithm is similar to the MRU-LRU algorithm, except that the most/least recently used task is replaced with the most/least frequently used task [8], [9].
- ***Dynamic-computing-frequency reactive service with no cache*** (**DFNC**): This algorithm provides reactive service to the requested task, where the mobile device first downloads the input data from the MEC server and then computes it to obtain the output data.
- ***Dynamic-computing-frequency reactive service with cache*** (**DFC**): This algorithm provides reactive service to the requested task, with the option of caching the input and output data into the limited capacity.

It is important to note that the DFC, DFNC, and PTDFC algorithms are all implemented with the SAC algorithm, and as a result, we refer to them as '*SAC-enabled algorithms*' in the following analysis.

#### B. Data Simulation

In this study, the training and testing data were generated through a simulation process involving the creation of a Markov chain from a set of tasks $\mathcal{F}$. The transit probability of a task $i$ to another randomly selected task $j \in \mathcal{F} \backslash i$ was established as the maximum transition probability, $p_{i,j} = p_{\max}$. For other tasks $k \in \mathcal{F} \backslash j$, the probability $p_{i,k}$ was calculated as $(1 - p_{i,j}) \frac{|p'_{i,k}|}{\sum_{f \in \mathcal{F} \backslash j} |p'_{i,f}|}$, where $p'_{i,k}$ or $p'_{i,f}$ were randomly sampled from a uniform distribution. The resulting Markov chain represented the request popularity and transition preferences of $F$ tasks. Subsequently, $10^6$ requested tasks were sampled using a frame-by-frame method. To account for the slow movement of mobile devices, the SNR of the communication channel was dynamically changed every 300 epochs, with four possible values: $0.5\,\text{dB}$, $1\,\text{dB}$, $2\,\text{dB}$, and $3\,\text{dB}$. The transition between

different SNRs was randomized with equal probabilities. The simulation was conducted using default configurations, which included $M = 8$, $F = 4$, a maximum transition probability of 0.7, $\lambda = 1$, $I_f$ around 16000 bits with random offset, $O_f$ around 30000 bits with random offset, $w = 800$ cycles/bit, $\tau = 0.02$ seconds, $f_D = 1.7 \times 10^8$ cycles/s, $\mu = 10^{-19}$, and $C = 40000$ bits.

## C. Implementation

For the purposes of training and stabilization, the SAC action $\mathbf{a}_t$ and system state $\mathbf{x}_t$ are normalized to fall within the range of $[-1, 1]$. Implementation of the system is accomplished through the use of Python and PyTorch. Training and testing processes are executed on a computer with a TITAN RTX GPU, utilizing a batch size of 256, a discount factor of $\gamma = 0.99$, automatic entropy temperature $\alpha$ tuning [29], a hidden-layer size of 256, one model update per step, one target update per 1000 steps, and a replay buffer size of $1 \times 10^7$. The testing process is executed 10 epochs after every 10 training epochs, and the training and testing processes are halted when the reward and loss have converged.



Fig. 2: Training reward of PTDFC, DFC, DFNC, MFU-LFU, MRU-LRU algorithms when the SNR values are dynamically changed every 300 epochs.

## D. Convergence Analysis

We present the training convergence results of three SAC-based algorithms, PTDFC, DFC, DFNC, and two heuristic algorithms, MFU-LFU and MRU-LRU in Fig. 2. The curves plot the reward versus epochs under different SNR conditions for these algorithms. It is important to note that the MFU-LFU and MRU-LRU algorithms are heuristic in nature, lacking parameters for training. Despite this, we have included their reward outcomes in Fig. 2, aiming to provide a more comprehensive perspective on their relative performance compared to others. During the first 300 epochs, with SNR $= 1$ dB, the PTDFC, DFC, and DFNC algorithms commence with neural network parameters initialized randomly and achieve convergence in a substantial number of epochs (250, 135, and 20 epochs, respectively). PTDFC requires more training epochs to converge than DFC and DFNC simply because of its larger action space. Starting from epoch 300, the SNR value is increased to $2$ dB, and the three SAC-based algorithms converge again in less than 13 epochs using the pre-trained model from the previous epochs. This finding demonstrates the remarkable generalization ability of SAC-based algorithms
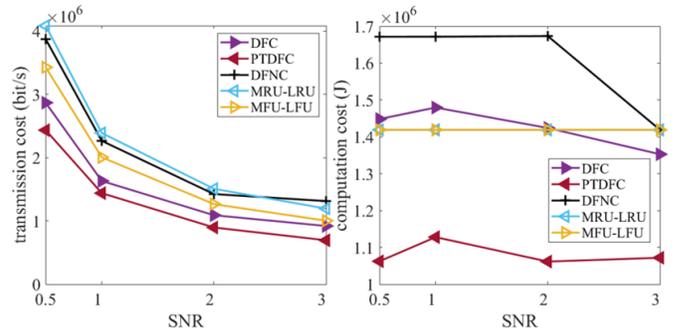


Fig. 3: The system performance for the proposed PTDFC algorithm and the baselines with the configuration stated in Section. V-B and V-C. (left) Transmission cost vs. SNR. (right) Computation cost vs. SNR.

to handle SNR change cases. These SAC-based algorithms can get fine-tuned and converged again within a few epochs (around 10). The quick convergence ability is also validated at epochs 600 and 900 when the SNR changes to $0.5$ dB and $3$ dB, respectively. We also noticed that there are significant discrepancies in the convergence time between the two SNR $= 1$ dB stages. In the later SNR $= 1$ dB stage (epoch 1200-1500), the PTDFC achieves convergence in approximately 10 epochs after the SNR transition. This can be attributed to the solid foundation of well-trained parameters established during the preceding SNR $= 3$ dB stage. This trend reaffirms the system's adeptness in rapidly adapting to environmental SNR changes.

The optimization problem at hand involves both linear and nonlinear objective functions, constraints, and involves binary variables. Consequently, it is classified as an Integer Nonlinear Programming (INLP) problem and is notoriously difficult to solve. Traditional optimization algorithms for INLP (such as Branch and Bound) are not suitable for this problem due to their exponential convergence time and the assumption of global knowledge of the environment and its dynamics. Moreover, in the event of a change in the environment, such as a variation in channel SNR, it takes a considerable amount of time to solve the problem and achieve convergence again. Classical machine learning algorithms, exceptionally standard reinforcement learning, also face challenges in scaling with the large dimension of the variable set, requiring an excessively large network size and convergence time to solve the problem.

## E. Numerical Results

The system performance of the proposed PTDFC algorithm and the baselines (DFC, DFNC, MRU-LRU, MFU-LFU) in terms of transmission bandwidth cost and computation cost for different channel SNRs is presented in Fig. 3. The results show that the PTDFC algorithm achieves the lowest cost for transmission bandwidth for various channel SNRs, followed by the DFC, MFU-LFU, DFNC, and MRU-LRU algorithms. In terms of computation cost, the top-performing algorithms are PTDFC, MFU-LFU/MRU-LRU, DFC, and DFNC, respectively. Overall, the PTDFC algorithm achieves a reduction of around $5 \times 10^5$ bits/s in transmission cost and $3 \times 10^5$ J in
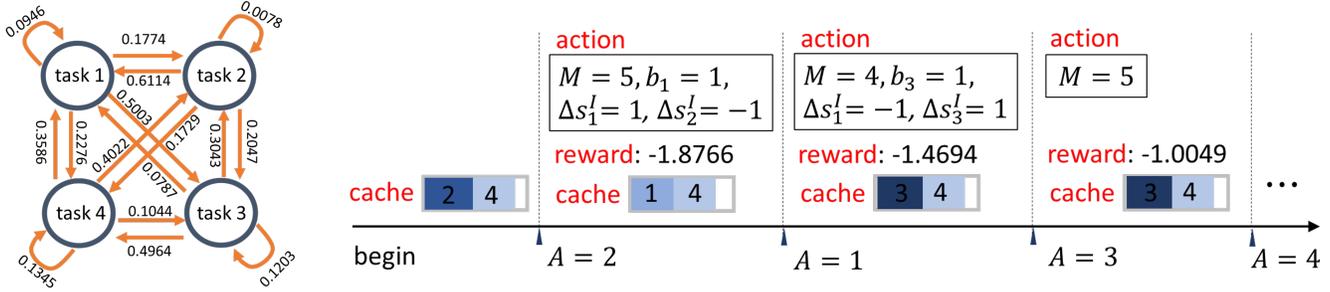
Fig. 4: A qualitative example for the joint optimization of 4 tasks using SAC-based PTDFC algorithm. (left) Visualization of the Markov transition probability among 4 tasks. (right) The requested task, cache state, action, and reward for the first four time slots of the proposed SAC system. If no action is mentioned, it defaults to no change with a value of 0.

computation cost for every SNR condition, compared to the second-best algorithm. It is also observed that all algorithms take less transmission bandwidth for the requested task as the SNR value increases, indicating that a higher SNR results in better channel quality.

### F. Qualitative Results Analysis

Fig. 4 provides an example of the status and action of four requests when deploying the PTDFC algorithm. The figure visualizes the requested task, cache state, action, and reward of each time slot to show the joint computing, pushing, and caching optimization of the four tasks. In this example, at $t_0$, the mobile device requests task 1 from the MEC server, which has empty cache content. The system then makes reactive transmission and computing for task 1 with five cores and pushes the input data of unrequested task 3, followed by caching the input data of tasks 1 and 3. At $t_1$, the requested task is task 3, and the system makes the reactive computing of the cached task 3 with four cores and pushes the input data of task 4. Then, the system replaces the cache of task 1 with the input data for task 4. Similarly, at $t_2$, the requested task is task 4, and the system makes the reactive computing of the cached task 4 with five cores and pushes the input data of task 2. Finally, the system removes the cache for task 3 and caches the input data for task 2. The example illustrates that the SAC-based PTDFC system is capable of predicting the user's future requests using deep networks and pushing or caching the appropriate content to enhance system performance.

### G. Tuning Analysis

In this section, we investigate the impact of several crucial parameters on the performance of the proposed PTDFC algorithm. These parameters include the cache size ($C$), number of computing cores ($M$), number of tasks ($F$), maximum transition probabilities, base computing frequency ($f_D$), task input size ($I_f$), task output size ($O_f$), tolerable service delays ($\tau$), and cost weights ($\lambda$). To analyze the effects of each parameter, we hold the other parameters constant and observe the resulting changes in performance. The default values for these parameters are specified in Section V-B, and we maintain a fixed channel SNR value of 1 to isolate the effects of parameter tuning.
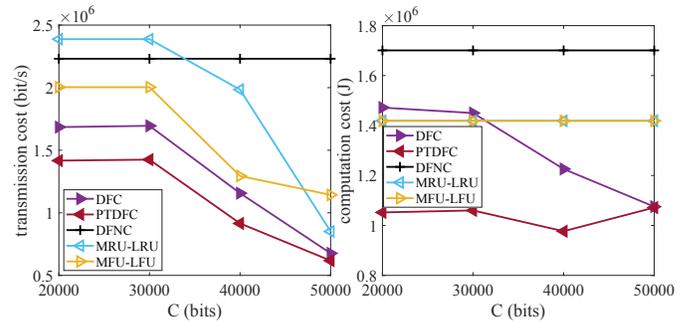


Fig. 5: Impact of varying the cache size $C$ when using the default configuration and fixed SNR. (left) Transmission cost vs. $C$. (right) Computation cost vs. $C$.

*1) Different Cache Size $C$:* In Fig. 5, we have presented the averaged transmission and computation costs of three SAC-enabled algorithms, DFNC, DFC, and PTDFC, as well as two heuristic algorithms, MRU-LRU and MFU-LFU, under different cache sizes $C$. It is worth noting that the DFNC algorithm is not affected by changes in cache size, as it only provides reactive service without caching. The other algorithms show a decrease in transmission costs as the cache size is increased, due to the availability of more locally cached input data. Moreover, our proposed PTDFC algorithm consistently achieves lower transmission and computation costs than the other algorithms, thanks to its ability to dynamically adjust the cache via proactive transmission. With a very large cache size, (e.g., $C = 50000$ bits), the performance of PTDFC and DFC is similar because the cache is large enough to store all input data and there is no need for a proactive transmission. Furthermore, we have observed a consistent overlapping trend in the computation costs of MRU-LRU and MFU-LFU across various configurations. This overlapping behavior can be attributed to our design choice in both algorithms, wherein solely the input data is cached. Consequently, the performance of MRU-LRU and MFU-LFU in terms of computation cost tends to align, as governed by Eq. (5), when $S_{A(t)}^O(t) = 0$.

*2) Different Number of Computation Cores $M$:* Fig. 6 illustrates the performance of five algorithms, namely DFNC, DFC, PTDFC, MRU-LRU, and MFU-LFU, under different
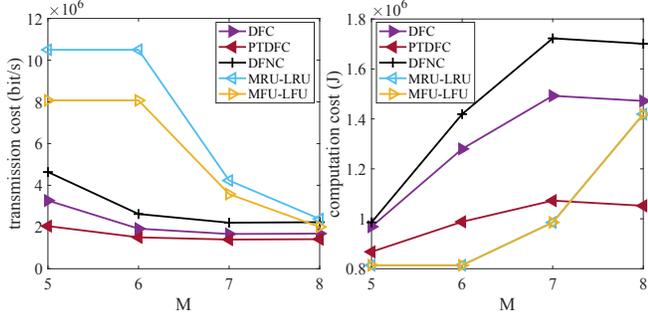
Fig. 6: Impact of varying the number of computing cores $M$ when using the default configuration and fixed SNR. (left) Transmission cost vs. $M$. (right) Computation cost vs. $M$.



Fig. 7: Impact of varying the task number $F$ when using the default configuration and fixed SNR. (left) Transmission cost vs. $F$. (right) Computation cost vs. $F$.

numbers of computation cores $M$. As the number of computing cores increases, the transmission cost of all five algorithms decreases, while their computation cost increases correspondingly. This is because all the reactive tasks are time-sensitive, and the system tends to utilize more computing cores for computing to reduce the computing time and leave more time for reactive transmission, which would effectively reduce the transmission cost and total cost. The proposed PTDFC algorithm consistently achieves a low transmission cost and computation cost by selecting an appropriate computing core number for the required task to achieve a better reward or a smaller cost. In contrast, the heuristic algorithms MRU-LRU and MFU-LFU have high transmission costs for small $M$ and significant computing costs for large $M$, as their computing frequency is linearly adjusted with the increase of computing cores.

*3) Different Number of Tasks $F$:* In Fig. 7, we present the performance of five algorithms under various numbers of tasks in the request set. As the number of tasks increases, we observe a slight increase in the transmission and computation costs for the DFC and PTDFC algorithms, possibly because only a small portion of all tasks can be cached or proactively transmitted, while the rest has to be reactively served, leading to higher costs. However, it is worth noting that the proposed PTDFC algorithm consistently outperforms all other algorithms across different numbers of tasks. This is due to its ability to dynamically select the best task in the task set for caching and proactive transmission, thereby minimizing costs associated with reactive service.

*4) Different Maximum Transition Probabilities:* We investigated the impact of the maximum transition probabilities $p_{\max}$ on the Markov chain simulation, which is a crucial parameter. Accordingly, we conducted experiments by varying the $p_{\max}$ value and evaluated the performance of five algorithms under different $p_{\max}$ values. The results, shown in Fig. 8, indicate that a higher $p_{\max}$ value leads to an easier prediction of future service requests based on the current request, resulting in the corresponding proactive push operation. Consequently, the PTDFC algorithm, equipped with proactive transmission, demonstrated significant reductions in transmission and computation costs compared to the other algorithms, which do
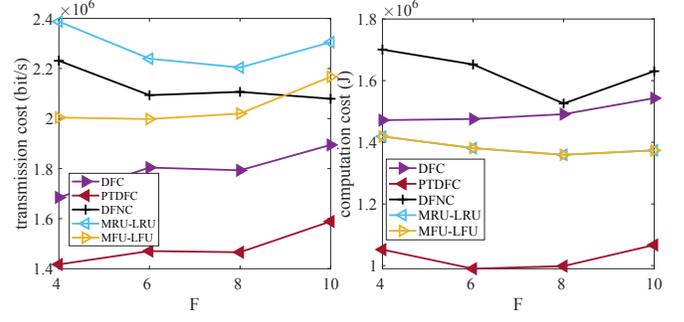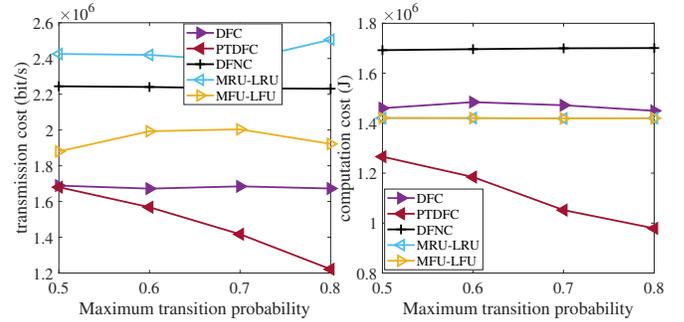


Fig. 8: Impact of varying the maximum transition probability in Markov chain when using the default configuration and fixed SNR. (left) Transmission cost vs. maximum transition probability. (right) Computation cost vs. maximum transition probability.

not consider proactive transmission. It is worth noting that the other algorithms exhibited little to no variation in their performance across different $p_{\max}$ values, as they do not rely on the push operation.

*5) Different Base Computing Frequency $f_D$:* To study the impact of computing frequency on system performance, we conducted an evaluation of five algorithms with varying $f_D$ and presented our results in Fig. 9. As a general rule, increasing the computing frequency is equivalent to having more computing cores with a fixed value of $f_D$. Therefore, the trends observed in the results of Fig. 9 are similar to those previously reported in Fig. 6. Notably, the PTDFC algorithm consistently exhibits the lowest overall cost (i.e., the sum of transmission and computation costs) across all $f_D$ configurations.

*6) Different Task Input Size $I_f$:* The default input data size for all tasks is 16000 bits. To investigate the impact of input data size on the system performance, we conducted experiments by varying the input data size for four tasks and evaluated the performance of five algorithms under different $I_f$. The results are presented in Fig. 10, where it can be observed that an increase in the input data size leads to an increase in both transmission and computation costs for all algorithms. This can be attributed to the fact that larger input
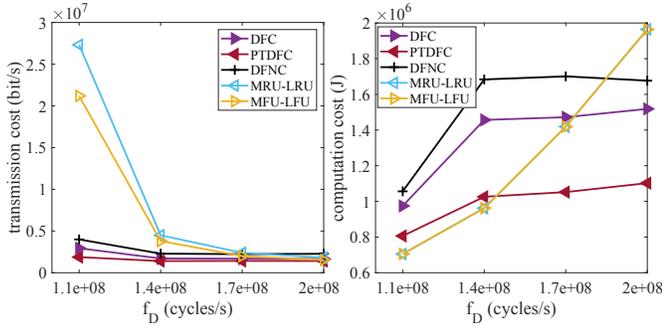
Fig. 9: Impact of varying the base computing frequency $f_D$ when using the default configuration and fixed SNR. (left) Transmission cost vs. $f_D$. (right) Computation cost vs. $f_D$.
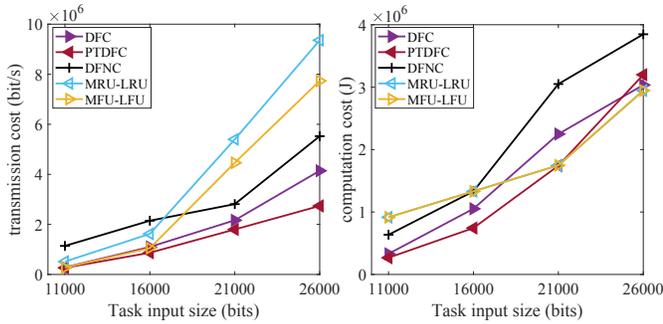


Fig. 11: Impact of varying the output data size $O_f$ for 4 tasks when using the default configuration and fixed SNR. (left) Transmission cost vs. $O_f$. (right) Computation cost vs. $O_f$.



Fig. 10: Impact of varying the input data size $I_f$ for 4 tasks when using the default configuration and fixed SNR. (left) Transmission cost vs. $I_f$. (right) Computation cost vs. $I_f$.

data requires more bandwidth and computation, whether in reactive or proactive cases, as indicated by Eq. (3), (5), (6). On the other hand, when the input data size is relatively small (e.g., 11000 bits), the PTDFC and DFC algorithms exhibit similar performance, as the optimal policy for both algorithms is to cache as much input data as possible and provide full reactive service only for non-cached tasks. However, in general, the PTDFC algorithm consistently outperforms the other algorithms, achieving the smallest total cost across all input data size configurations.

*7) Different Task Output Size $O_f$:* After analyzing the impact of input data size, we further examined the influence of output data size on the system performance. For this purpose, we altered the output data size of four tasks and evaluated the performance of five algorithms under different $O_f$. The default output data size for all tasks was set to 30000 bits. As depicted in Fig. 11, the transmission cost and computation cost for the DFNC, MRU-LRU, and MFU-LFU algorithms remained unchanged, as the output data size did not affect the calculation of the two costs and the cache update mechanism. However, for the DFC algorithm, the transmission cost fluctuated around a constant level, and the computation cost increased as the output data size increased from 15000 bits to 30000 bits. This behavior can be attributed to the algorithm's prioritization of caching input data, which ensures a low transmission cost.
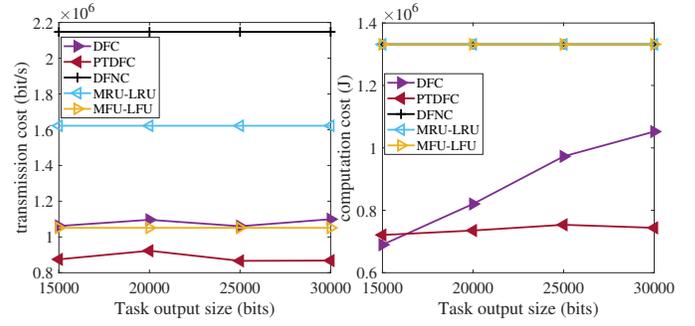
On the other hand, due to limited cache size, lower priority, and increased output data size, only a small fraction of tasks had the opportunity to cache their output data, resulting in additional computation for reactive service. In contrast, the PTDFC algorithm's joint computing, pushing, and caching design demonstrated robustness to variations in $O_f$, as evidenced by the flat transmission and computation cost curves in Fig. 11.

*8) Different Tolerable Service Delays $\tau$:* The tolerable service delay is a critical parameter that significantly influences the transmission and computation costs of the system. To evaluate their impact, we tested the performance of five algorithms under varying values of $\tau$ and present our findings in Fig. 12. As we increase $\tau$ from $0.012\,\mathrm{s}$ to $0.024\,\mathrm{s}$, we observe a corresponding decline in the transmission and computation costs for most algorithms. This is because the larger $\tau$ provides more time for transmitting the input data of the requested task, reducing the bandwidth cost as per Eq. (3). Similarly, additional processing time is given to the computation step for acquiring the output data, relaxing the requirement of the computing frequency, and leading to a lower computation cost as per Eq. (5). Notably, the PTDFC algorithm outperforms the other algorithms by achieving the lowest transmission and computation cost under all $\tau$ values. This is attributed to its ability to design an optimal policy for joint computing, pushing, and caching through deep reinforcement learning. However, the transmission cost of all five algorithms converges at larger $\tau$, and the benefits of the PTDFC algorithm are mitigated as a lower computing frequency (one core) is employed for all algorithms, resulting in comparable transmission costs.

*9) Different Cost Weights $\lambda$:* To guide policy learning for the trade-off between transmission cost and computation cost, the default cost weight of 1 is used for designing the reward function in Eq. (23). To investigate the impact of this parameter on the performance of the algorithms, we evaluated the performance of five algorithms under different $\lambda$ values and present the results in Fig. 13. Notably, the SAC-enabled algorithms show a clear trade-off between the two costs, as increasing $\lambda$ leads to a decrease in computation cost and a corresponding increase in transmission cost, while the heuristic algorithms MRU-LRU and MFU-LFU have completely flat cost curves. The PTDFC algorithm consistently achieves the
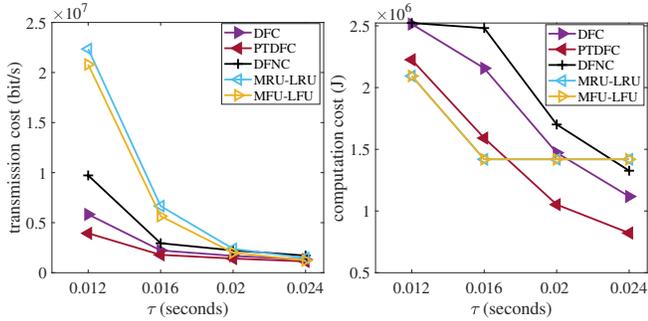
Fig. 12: Impact of varying the maximum tolerable service latency $\tau$ when using the default configuration and fixed SNR. (left) Transmission cost vs. $\tau$. (right) Computation cost vs. $\tau$.
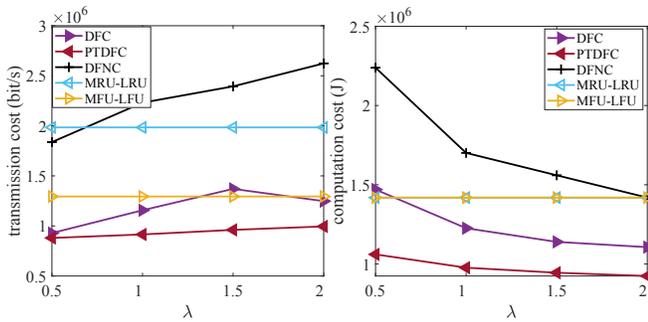


Fig. 13: Impact of varying cost weight $\lambda$ when using the default configuration and fixed SNR. (left) Transmission cost vs. $\lambda$. (right) Computation cost vs. $\lambda$.

best performance under different $\lambda$ values, which can be attributed to its optimal computing, pushing, and caching policy design through deep reinforcement learning.

### H. Complexity Analysis

The computational complexity of the proposed PTDFC algorithm largely relies on the number and structure of neural networks in SAC system [35]. At the training stage, Algorithm 1 incorporates the parameter updating of three neural networks: $Q_\theta$, $Q_{\bar{\theta}}$ (the actors), and $\pi_\phi$ (the critic). Therefore, the computation of the complexity of Algorithm 1 is:

$$2 \times \sum_{j=0}^{J-1} n_j^Q n_{j+1}^Q + \times \sum_{k=0}^{K-1} n_k^\pi n_{k+1}^\pi$$
$$= O\left( \sum_{j=0}^{J-1} n_j^Q n_{j+1}^Q + \sum_{k=0}^{K-1} n_k^\pi n_{k+1}^\pi \right) \quad (24)$$

where $J$ denotes the number of fully connected layers for the $Q_\theta$ and $Q_{\bar{\theta}}$ networks (having identical structure), and $K$ denotes that for $\pi_\phi$ network. $n_j^Q$ and $n_k^\pi$ represent the number of neurons at the $j$-th layer of $Q_\theta$ or $Q_{\bar{\theta}}$ networks and the $k$-th layer of $\pi_\phi$ network. $j = 0$ and $k = 0$ represent the input layers.

At the testing stage, Algorithm 1 only needs to execute the trained $Q_\theta$, $Q_{\bar{\theta}}$ networks, so the computation complexity is reduced to $O\left( \sum_{j=0}^{J-1} n_j^Q n_{j+1}^Q \right)$. In our system, $J = 3$, $K = 4$, $N_j^Q = 22, 256, 256, 1$ for $j = 0, 1, 2, 3$, given the number of tasks $F = 4$.

## VI. CONCLUSION

In this paper, we explore joint optimization of computing, pushing, and caching in MEC networks to further improve user-perceived quality of experience. We formulate the joint-design problem as an infinite-horizon discounted-cost Markov decision process, which allows us to optimize the total quantity of transmitted data and the total computation cost for the mobile user. To solve this problem, we propose a framework based on SAC learning that dynamically orchestrates the three functions. The framework is featured with embedded deep networks that implicitly predict user future requests and a design for action quantization and correction that enables SAC to work for this problem. In simulations using a single-user single-server MEC network, our proposed framework effectively reduces both transmission load and computing cost and outperforms baseline algorithms across various parameters.

## REFERENCES

[1] Y. Sun, Z. Chen, M. Tao, and H. Liu, "Communications, caching, and computing for mobile virtual reality: Modeling and tradeoff," *IEEE Trans. Commun.*, vol. 67, no. 11, pp. 7573–7586, Nov. 2019.

[2] X. Gao, G. Xing, S. Roy, and H. Liu, "Experiments with mmwave automotive radar test-bed," in *2019 53rd Asilomar Conference on Signals, Systems, and Computers*, 2019, pp. 1–6.

[3] X. Gao, S. Roy, G. Xing, and S. Jin, "Perception through 2d-mimo fmcw automotive radar under adverse weather," in *2021 IEEE International Conference on Autonomous Systems (ICAS)*, 2021, pp. 1–5.

[4] X. Gao, S. Ding, K. Vanas, D. R. Harshavardhan, and H. Soder-lund, "Deformable radar polygon: A lightweight and predictable occupancy representation for short-range collision avoidance," *arXiv preprint arXiv:2203.01442*, 2022.

[5] U. Cisco, "Cisco annual internet report (2018-2023)," *White Paper*, 2020. [Online]. Available: https://www.cisco.com/c/en/us/solutions/collateral/executive-perspectives/annual-internet-report/white-paper-c11-741490.html

[6] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Transactions on information theory*, vol. 60, no. 5, pp. 2856–2867, 2014.

[7] Y. Sun, Z. Chen, and H. Liu, "Delay analysis and optimization in cache-enabled multi-cell cooperative networks," in *IEEE Global Communications Conference (GLOBECOM)*, 2016, pp. 1–7.

[8] J. Wang, "A survey of web caching schemes for the internet," *SIGCOMM Comput. Commun. Rev.*, vol. 29, no. 5, p. 36–46, oct 1999. [Online]. Available: https://doi.org/10.1145/505696.505701

[9] Y. Sun, Y. Cui, and H. Liu, "Joint pushing and caching for bandwidth utilization maximization in wireless networks," *IEEE Transactions on Communications*, vol. 67, no. 1, pp. 391–404, 2019.

[10] W. Chen and H. V. Poor, "Content pushing with request delay information," *IEEE Transactions on Communications*, vol. 65, no. 3, pp. 1146–1161, 2017.

[11] Y. Lu, W. Chen, and H. V. Poor, "Coded joint pushing and caching with asynchronous user requests," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 8, pp. 1843–1856, 2018.

[12] M. Gregori, J. Gómez-Vilardebó, J. Matamoros, and D. Gündüz, "Wireless content caching for small cell and d2d networks," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 5, pp. 1222–1234, 2016.

[13] X. Yang, Z. Fei, J. Zheng, N. Zhang, and A. Anpalagan, "Joint multi-user computation offloading and data caching for hybrid mobile cloud/edge computing," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 11, pp. 11 018–11 030, 2019.

[14] M. Chen, Y. Hao, L. Hu, M. S. Hossain, and A. Ghoneim, "Edge-cocaco: Toward joint optimization of computation, caching, and communication on edge cloud," *IEEE Wireless Communications*, vol. 25, no. 3, pp. 21–27, 2018.

[15] Y. Hao, M. Chen, L. Hu, M. S. Hossain, and A. Ghoneim, "Energy efficient task caching and offloading for mobile edge computing," *IEEE Access*, vol. 6, pp. 11 365–11 373, 2018.

[16] L. Zhang, Y. Sun, Z. Chen, and S. Roy, "Communications-caching-computing resource allocation for bidirectional data computation in mobile edge networks," *IEEE Transactions on Communications*, vol. 69, no. 3, pp. 1496–1509, 2021.

[17] Y. Sun, Z. Chen, M. Tao, and H. Liu, "Bandwidth gain from mobile edge computing and caching in wireless multicast systems," *IEEE Transactions on Wireless Communications*, vol. 19, no. 6, pp. 3992–4007, 2020.

[18] X. Gao, G. Xing, S. Roy, and H. Liu, "Ramp-cnn: A novel neural network for enhanced automotive radar object recognition," *IEEE Sensors Journal*, vol. 21, no. 4, pp. 5119–5132, 2021.

[19] X. Gao, H. Liu, S. Roy, G. Xing, A. Alansari, and Y. Luo, "Learning to detect open carry and concealed object with 77 ghz radar," *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 4, pp. 791–803, 2022.

[20] Y. Qian, R. Wang, J. Wu, B. Tan, and H. Ren, "Reinforcement learning-based optimal computing and caching in mobile edge network," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 10, pp. 2343–2355, 2020.

[21] Y. Wei, F. R. Yu, M. Song, and Z. Han, "Joint optimization of caching, computing, and radio resources for fog-enabled iot using natural actor–critic deep reinforcement learning," *IEEE Internet of Things Journal*, vol. 6, no. 2, pp. 2061–2073, 2019.

[22] Z. Ning, K. Zhang, X. Wang, M. S. Obaidat, L. Guo, X. Hu, B. Hu, Y. Guo, B. Sadoun, and R. Y. K. Kwok, "Joint computing and caching in 5g-envisioned internet of vehicles: A deep reinforcement learning-based traffic control system," *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 8, pp. 5201–5212, 2021.

[23] L. Huang, X. Feng, A. Feng, Y. Huang, and L. P. Qian, "Distributed deep learning-based offloading for mobile edge computing networks," *Mob. Netw. Appl.*, vol. 27, no. 3, p. 1123–1130, jun 2022. [Online]. Available: https://doi.org/10.1007/s11036-018-1177-x

[24] R. Zhao, X. Wang, J. Xia, and L. Fan, "Deep reinforcement learning based mobile edge computing for intelligent internet of things," *Physical Communication*, vol. 43, p. 101184, 2020. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1874490720302615

[25] L. Huang, S. Bi, and Y.-J. A. Zhang, "Deep reinforcement learning for online computation offloading in wireless powered mobile-edge computing networks," *IEEE Transactions on Mobile Computing*, vol. 19, no. 11, pp. 2581–2593, 2020.

[26] X. Gao, S. Roy, and G. Xing, "Mimo-sar: A hierarchical high-resolution imaging algorithm for mmwave fmcw radar in autonomous driving," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 8, pp. 7322–7334, 2021.

[27] X. Gao, S. Roy, and L. Zhang, "Static background removal in vehicular radar: Filtering in azimuth-elevation-doppler domain," *arXiv preprint arXiv:2307.01444*, 2023.

[28] P. Dai, F. Song, K. Liu, Y. Dai, P. Zhou, and S. Guo, "Edge intelligence for adaptive multimedia streaming in heterogeneous internet of vehicles," *IEEE Transactions on Mobile Computing*, vol. 22, no. 3, pp. 1464–1478, 2023.

[29] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel *et al.*, "Soft actor-critic algorithms and applications," *arXiv preprint arXiv:1812.05905*, 2018.

[30] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *nature*, vol. 518, no. 7540, pp. 529–533, 2015.

[31] H. Yin, L. Zhang, and S. Roy, "Multiplexing urllc traffic within embb services in 5g nr: Fair scheduling," *IEEE Transactions on Communications*, vol. 69, no. 2, pp. 1080–1093, 2021.

[32] H. Yin, P. Liu, K. Liu, L. Cao, L. Zhang, Y. Gao, and X. Hei, "Ns3-ai: Fostering artificial intelligence algorithms for networking research," in *Proceedings of the 2020 Workshop on Ns-3*, ser. WNS3 '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 57–64. [Online]. Available: https://doi.org/10.1145/3389400.3389404

[33] L. Zhang, H. Yin, Z. Zhou, S. Roy, and Y. Sun, "Enhancing wifi multiple access performance with federated deep reinforcement learning," in *2020 IEEE 92nd Vehicular Technology Conference (VTC2020-Fall)*, 2020, pp. 1–6.

[34] K. Psounis, A. Zhu, B. Prabhakar, and R. Motwani, "Modeling correlations in web traces and implications for designing replacement policies," *Computer Networks*, vol. 45, no. 4, pp. 379–398, 2004.

[35] F. Zhang, G. Han, L. Liu, M. Martinez-Garcia, and Y. Peng, "Deep reinforcement learning based cooperative partial task offloading and resource allocation for iiot applications," *IEEE Transactions on Network Science and Engineering*, pp. 1–1, 2022.