Soft Actor-Critic Learning-Based Joint Computing, Pushing, and Caching Framework in MEC Networks

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

Abstract—To support future 6G mobile applications, the mobile edge computing (MEC) network needs to be jointly optimized for computing, pushing, and caching to reduce transmission load and computation cost. To achieve this, we propose a framework based on deep reinforcement learning that enables the dynamic orchestration of these three activities for the MEC network. The framework can implicitly predict user future requests using deep networks and push or cache the appropriate content to enhance performance. To address the curse of dimensionality resulting from considering three activities collectively, we adopt the soft actor-critic reinforcement learning in continuous space and design the action quantization and correction specifically to fit the discrete optimization problem. We conduct simulations in a single-user single-server MEC network setting and demonstrate that the proposed framework effectively decreases both transmission load and computing cost under various configurations of cache size and tolerable service delay.

I. INTRODUCTION

Recent advancements in smart mobile devices have enabled various emerging applications, such as virtual reality (VR) and augmented reality (AR), which require ultra-high communication and computation capabilities in low latency. To minimize these costs while ensuring a high-quality user experience, the MEC network is a promising solution that can push caching and computing resources to access points, base stations, and even mobile devices at the wireless network edge.

Caching can improve bandwidth utilization by placing frequently accessed content closer to users for future use,

X. Gao is with the University of Washington, Seattle, 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) Y. Sun is also with the Future Network of Intelligent Institute (FNii), the Chinese University of Hong Kong (Shenzhen), Shenzhen 518172, China. X. Xu and P. Zhang are 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, pzhang@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 affiliated with the Department of Broadband Communication, Peng Cheng Laboratory, Shenzhen 518172, China. S. Cui is also affiliated with the Department of Broadband Communication, Shenzhen 518000, China (email: shuguangcui@cuhk.edu.cn).

which is particularly useful due to the high degree of asynchronous content reuse in mobile traffic. Caching policies can be categorized into two types: *static caching* and *dynamic caching*. Static caching policies are generally based on content popularity distribution and involve cache states that remain unchanged over a relatively long period [1]. In contrast, dynamic caching policies involve content placement updates based on instantaneous user request information, such as the least recently used (LRU) and least frequently used (LFU) policy [2].

A joint pushing and caching design can improve system performance by proactively transmitting content during lowtraffic times to satisfy future user demands. Various joint pushing and caching designs exist that aim to maximize bandwidth utilization [3], effective throughput [4], minimize traffic load [5], or reduce transmit energy consumption [6]. However, these policies only consider content delivery and do not account for computation, therefore cannot be directly applied to modern mobile traffic services, such as VR delivery.

To effectively serve mobile traffic, previous designs have considered the joint utilization of cache and computing resources at MEC servers to minimize transmission latency [7], [8] or energy consumption [9]. Some designs also aim to minimize transmission data [10]. However, these designs only consider static caching and do not allow for pushing.

To address the aforementioned issues, we propose a joint computing, pushing, and caching policy optimization approach and validate it in a single-user single-server MEC network. Our approach involves the following steps: (1) We formulate the joint optimization problem as an infinite-horizon discounted Markov decision process, where the aim is to minimize both computation cost and transmission dataload. (2) We use the soft actor-critic (SAC) reinforcement learning (RL) algorithm [11] to quickly and stably obtain dynamic computing, pushing, and caching policies. Unlike the classic deep Q-learning algorithm, which requires a Q-network with output nodes for all potential actions, SAC learns Q-functions with few parameters, addressing the curse of dimensionality. We designed an action quantization and correction mechanism to enable SAC, which operates in continuous space, to meet our discrete optimization requirements. (3) We present simulation results with various system parameters to demonstrate the effectiveness of our proposed algorithm.

II. SYSTEM MODEL

The MEC network we consider comprises a server and a mobile device with caching and computing capabilities, as

The work was supported in part by NSFC with Grant No. 62293482, the Basic Research Project No. HZQB-KCZYZ-2021067 of Hetao Shenzhen-HK S&T Cooperation Zone, NSFC with Grant No. 62301471, The Major Key Project of PCL Department of Broadband Communiation, the National Key R&D Program of China with grant No. 2018YFB1800800, the Shenzhen Outstanding Talents Training Fund 202002, the Guangdong Research Projects No. 2017ZT07X152 and No. 2019CX01X104, the Guangdong Provincial Key Laboratory of Future Networks of Intelligence (Grant No. 2022B1212010001), the Shenzhen Key Laboratory of Big Data and Artificial Intelligence (Grant No. ZDSYS201707251409055), and the Key Area R&D Program of Guangdong Province with grant No. 2018B030338001.



Fig. 1: Illustration of MEC network with single MEC server and single mobile device.

illustrated in Fig. 1. The MEC server's cache size is large enough to proactively store all input and output data related to tasks requested by the mobile device. In contrast, the mobile device's cache size is limited to C bits. The mobile device is equipped with multi-core computing capabilities, with each core operating at a frequency of f_D cycles per second. We assume that the mobile device has M computing cores. The system operates over an infinite time horizon, and time is slotted in intervals of τ seconds, with time slots indexed as $t = 0, 1, 2, \cdots$. At the start of each time slot, the mobile device submits a task request that is assumed to be delay-intolerant and must be served before the slot ends.

A. Task Model

We consider a set of F tasks that can be requested by the mobile device and denote this set as $\mathcal{F} \stackrel{\Delta}{=} \{1, 2, \dots, f, \dots, F\}$. Specifically, each task $f \in \mathcal{F}$ is characterized by a 4-tuple $\{I_f \text{ (in bits)}, O_f \text{ (in bits)}, w_f \text{ (in cycles/bit)}, \tau \text{ (in seconds)}\}$. system state $\mathbf{X}(t) = (A(t), \mathbf{S}(t))$, the task request A(t) is Specifically, I_f represents the size of the remote input data that can be proactively generated from the Internet and cached. The output data size is represented by O_f . The parameters w_f and τ denote the required computation cycles per bit and the maximum allowable service latency.

B. System State

1) Request State: At the beginning of each time slot t, the mobile device generates one task request. Let $A(t) \in$ \mathcal{F} denote the request state of the mobile device, where A(t) = f represents that the mobile device requests task f. The cardinality of \mathcal{F} is F. Assume that A(t) evolves according to a first-order F-state Markov chain, denoted as $\{A(t): t = 0, 1, 2, ...\}$, which captures both task popularity and inter-task correlation of order one of the task demand process. Let $\Pr[A(t+1) = j | A(t) = i]$ denote the transition probability of going to state $j \in \mathcal{F}$ at time slot t + 1, given that the request state at time slot t is $i \in \mathcal{F}$ for the task demand process. Assume that $\{A(t)\}\$ is time-homogeneous. Denote with $\mathbf{Q} \triangleq (q_{i,j})_{i \in \mathcal{F}, j \in \mathcal{F}}$ the transition probability matrix of $\{A(t)\}$, where $q_{i,j} \triangleq \Pr[A(t+1) = j|A(t) = i]$. Furthermore, we restrict our attention to irreducible Markov chain and denote with $\mathbf{p} \triangleq (p_f)_{f \in \mathcal{F}}$ the limiting distribution of $\{A(t)\}$, where $p_f \triangleq \lim_{t\to\infty} \Pr[A(t) = f]$. Note 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 cache state of the input data for task f in the storage of the mobile device, where $S_f^I(t) = 1$ means that the input data for task f is cached in the mobile device, and $S_f^I(t) = 0$, otherwise. Let $S_f^O(t) \in$

 $\{0,1\}$ denote the cache state of the output data for task f in the storage of the mobile device, where $S_f^O(t) = 1$ means that the output data for task f is cached in the mobile device, and $S_{f}^{O}(t) = 0$, otherwise. Denote with C (in bits) the size of the cache space at the mobile device, and the cache size constraint is given by

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

Let $\mathbf{S}(t) \triangleq (S_f^I(t), S_f^O(t))_{f \in \mathcal{F}} \in \mathcal{S}$ denote the cache state of the mobile device at time slot t, where $S \triangleq$ $\{(S_f^I, S_f^O)_{f \in \mathcal{F}} \in \{0, 1\}^F \times \{0, 1\}^F : \sum_{f \in \mathcal{F}} I_f S_f^I + O_f S_f^O \leq C\} \text{ represents the cache state space of the mobile device. The cardinality of <math>\mathcal{S}$ is bounded by $\binom{F}{N_{\min}}$ and $\binom{F}{N_{\max}}$ from below and above, respectively, where $N_{\min} \triangleq \frac{C}{\max_{f \in \mathcal{F}} \{I_f, O_f\}}$, and $N_{\max} \triangleq \frac{C}{\min_{f \in \mathcal{F}} \{I_f, O_f\}}.$ 3) System State: At time slot t, the system state consists

of both system request state and system cache state, denoted as $\mathbf{X}(t) \triangleq (A(t), \mathbf{S}(t)) \in \mathcal{F} \times \mathcal{S}$, where $\mathcal{F} \times \mathcal{S}$ represents the system state space.

C. System Action

1) Reactive Computation Action: At time slot t, we denote with $B^{R}(t)$ and $E^{R}(t)$ the reactive transmission bandwidth cost and the reactive computation energy cost. Based on the served as below:

- If $S_{A(t)}^{O}(t) = 1$, the output of task A(t) can be directly obtained from the local cache without any transmission or computation. In this way, the delay is negligible, and the reactive computation energy or transmission cost is zero.
- If $S_{A(t)}^{I}(t) = 1$ and $S_{A(t)}^{O}(t) = 0$, the mobile device can directly compute the task based on the locally cached input data. Let $c_{R,f}(t) \in \{1, \dots, M\}$ denote the number of computation cores at the mobile device allocated for reactively processing task f at time slot t. Thus, we directly have $c_{R,f}(t) = 0$ for $\forall f \in \mathcal{F} \setminus A(t)$. In order to serve the requested task A(t) within τ , we assume that $\frac{I_f w_f}{\tau} \mathbf{1}(A(t) = f) \leq M f_D, \forall f \in \mathcal{F}$ and require $\frac{I_{A(t)}w_{A(t)}}{\tau} \leq c_{R,A(t)}(t)f_D$. The energy consumed for computing one cycle with frequency $c_{R,f}(t)f_D$ at the mobile device is $\mu c_{R,f}^2(t) f_D^2$, where μ is the effective switched capacitance related to the chip architecture indicating the power efficiency of CPU at the mobile device. 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 should first download the input data of the task A(t) from the MEC server, and then compute it locally. The required latency is given by $\frac{I_{A(t)}}{B^R(t)} + \frac{I_{A(t)}w_{A(t)}}{c_{R,A(t)}(t)f_D}$. In order to satisfy the latency constraint, i.e., $\frac{I_{A(t)}}{B^R(t)} + \frac{I_{A(t)}w_{A(t)}}{c_{R,A(t)}(t)f_D} \leq \tau$, the minimum reactive transmission cost $B^R(t)$ is given by

 $\frac{I_{A(t)}}{\tau - \frac{I_{A(t)}w_{A(t)}}{c_{R,A(t)}(t)f_D}}, \text{ and the reactive computation energy cost} \\ E^R(t) \text{ is given by } \mu c_{R,A(t)}^2(t) f_D^2 I_{A(t)} w_{A(t)}.$

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

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

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

$$B^{R}(t) = \left(1 - S^{I}_{A(t)}(t)\right) \left(1 - S^{O}_{A(t)}(t)\right) \frac{I_{A(t)}}{\tau - \frac{I_{A(t)}w_{A(t)}}{c_{R,A(t)}(t)f_{D}}},$$
(3)

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

$$E^{R}(t) = \left(1 - S^{O}_{A(t)}(t)\right) \mu c^{2}_{R,A(t)}(t) f^{2}_{D} I_{A(t)} w_{A(t)}.$$
 (4)

Denote with $\mathbf{c}_R \triangleq (c_{R,f})_{f \in \mathcal{F}} \in \Pi^R_C(\mathbf{X})$ the system reactive computation action, where $\Pi^R_C(\mathbf{X}) \triangleq$ $\left\{ (c_{R,f})_{f \in \mathcal{F}} \in \{0, 1, \cdots, M\}^F : (2) \right\}$ denotes the system reactive computation decision space under system state X. From (2), we can see that the cardinality of reactive computation action space is M+1.

2) Proactive Transmission or Pushing Action: Denote with $b_f(t) \in \{0,1\}$ the pushing decision of task $f \in \mathcal{F}$, where $b_f(t) = 1$ means that the remote input data of task f is pushed to the mobile device, and $b_f(t) = 0$, otherwise. Assume that by the end of the time slot, the pushed data are transmitted to the mobile device. In order to satisfy the latency constraint, we have $\frac{\sum_{f=1}^{F} I_f b_f(t)}{\tau} \leq B^P(t)$, where $B^P(t)$ denotes the proactive transmission bandwidth cost. Thus, the minimum proactive transmission cost is given by

$$B^{P}(t) = \frac{\sum_{f=1}^{F} I_{f} b_{f}(t)}{\tau}.$$
 (5)

In summary, denote with $\mathbf{b} \triangleq (b_f)_{f \in \mathcal{F}} \in \{0, 1\}^F$ the system pushing action, where 2^F represents the system pushing action space under system state X.

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

$$S_{f}^{I}(t+1) = S_{f}^{I}(t) + \Delta s_{f}^{I}(t), \tag{6}$$

$$S_{f}^{O}(t+1) = S_{f}^{O}(t) + \Delta S_{f}^{O}(t), \tag{7}$$

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 $\forall f \in \mathcal{F}$

$$-S_{f}^{I}(t) \leq \Delta s_{f}^{I}(t) \leq \min \left\{ b_{f}(t) + c_{R,f}(t), 1 - S_{f}^{I}(t) \right\}$$
(8)
$$S_{f}^{O}(t) \leq \Delta s_{f}^{O}(t) \leq \min \left\{ c_{f}(t), 1 - S_{f}^{O}(t) \right\}$$
(9)

$$-S_f(t) \le \Delta s_f(t) \le \min\left\{c_{R,f}(t), 1 - S_f(t)\right\}, \tag{9}$$

$$\sum_{f=1}^{r} I_f \left(S_f^f(t) + \Delta s_f^f(t) \right) + O_f \left(S_f^F(t) + \Delta s_f^F(t) \right) \le C, \quad (10)$$

where the left-hand side of Eq. (8) indicates that only when the input of task f is cached at the mobile device, can it be removed, the right-hand side of Eq. (8) indicates that only when the input of task f is proactively transmitted from the MEC server or is reactively transmitted, i.e., $b_f(t) = 1$ or $c_{R,f}(t) > 0$, and has not been cached before, can it be cached into the mobile device. The left-hand side of Eq. (9) indicates that only when the output of task f is cached at the mobile device, can it be removed, and the right-hand side of Eq. (9) indicates that only when the output of task f is reactively computed at the mobile device, i.e., $c_{R,f}(t) > 0$, and has not been cached before, can it be cached into the mobile device. Eq. (10) indicates that the updated cache state should satisfy the cache size constraint.

In summary, denote with $\Delta \mathbf{s} \triangleq \left(\Delta s_f^I, \Delta s_f^O\right)_{f \in \mathcal{F}}$ $\Pi_{\Delta s}(\mathbf{X}) \text{ the system cache update action, where } \Pi_{\Delta s}(\mathbf{X}) \triangleq \left\{ \left(\Delta s_f^I, \Delta s_f^O \right)_{f \in \mathcal{F}} \in \{-1, 0, 1\}^F \times \{-1, 0, 1\}^F : (8), (9), (10) \right\}$ denotes the system cache update action space under system state X.

4) System Action: At each time slot, the system action consists of the reactive computation action, proactive computation action, pushing action, and cache update action, denoted as $(\mathbf{c}_R, \mathbf{b}, \Delta \mathbf{s}) \in \Pi(\mathbf{X}), \text{ where } \Pi(\mathbf{X}) \triangleq \Pi_C^R(\mathbf{X}) \times \{0, 1\}^F \times$ $\Pi_{\Delta s}(\mathbf{X})$ denotes the system action space under system state Х.

D. System Cost

At time slot t, the system cost consists of the transmission bandwidth cost and the computation energy cost. In particular, the transmission bandwidth cost consists of both the reactive and proactive transmission costs, given by

$$B(t) = B^{R}(t) + B^{P}(t),$$
(11)

where $B^{R}(t)$ is given in Eq. (3) and $B^{P}(t)$ is given in Eq. (5). The computation energy cost is the reactive computation cost only, i.e.,

$$E(t) = E^R(t), \tag{12}$$

where $E^{R}(t)$ is given in Eq. (4). To balance the communication and computation cost, we choose the weighted sum $B(t) + \lambda E(t)$ as the system cost at time slot t.

III. PROBLEM FORMULATION

Given an observed system state 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 π is a mapping from system state X to system action $(\boldsymbol{c}_R, \boldsymbol{b}, \Delta \boldsymbol{s})$, i.e., $(\boldsymbol{c}_R, \boldsymbol{b}, \Delta \boldsymbol{s}) = \pi(\boldsymbol{X}) \in \Pi(\boldsymbol{X})$.

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

$$\phi(\pi) \triangleq \limsup_{T \to \infty} \sum_{t=0}^{T-1} \gamma^t \mathbb{E} \left[B(t) + \lambda E(t) \right], \qquad (13)$$

where the expectation is taken over the task request process.

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.

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

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

IV. SOFT ACTOR-CRITIC LEARNING

A. SAC System State and Action

The system state \mathbf{x} of SAC is designed to 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 $\mathbf{a} = (\bar{\mathbf{c}}_R, \bar{\mathbf{b}}, \Delta \bar{\mathbf{s}}) \in \bar{\Pi}(\mathbf{X}) \triangleq \bar{\Pi}_C^R(\mathbf{X}) \times [0, 1]^F \times \bar{\Pi}_{\Delta s}(\mathbf{X}).$

As $\bar{\mathbf{c}}_R \triangleq \{(\bar{c}_{R,f})_{f \in \mathcal{F}}\}$ must always equal zero for $f \in \mathcal{F} \setminus A(t)$, 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 [11]. 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

$$J(\pi) = \sum_{t=0}^{T} \mathbb{E}_{(\mathbf{x}_t, \mathbf{a}_t) \sim \rho_{\pi}} \left[r\left(\mathbf{x}_t, \mathbf{a}_t\right) + \alpha \mathcal{H}\left(\pi\left(\cdot \mid \mathbf{x}_t\right)\right) \right], \quad (14)$$

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

The SAC algorithm [11] is a policy iteration approach designed to solve the optimization problem in Eq. (14). It comprises two essential components: soft Q-function $Q_{\theta}(\mathbf{x}_t, \mathbf{a}_t)$, and policy $\pi_{\phi}(\mathbf{a}_t | \mathbf{x}_t)$. To deal with the large continuous domains, neural networks approximate these components, with the network parameters denoted by θ and ϕ . 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 [11], the update rules for θ and ϕ are provided below.

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

$$J_{Q}(\theta) = \mathbb{E}_{(\mathbf{x}_{t},\mathbf{a}_{t})\sim\mathcal{D}} \left[\frac{1}{2} \left(Q_{\theta} \left(\mathbf{x}_{t},\mathbf{a}_{t} \right) - \left(r\left(\mathbf{x}_{t},\mathbf{a}_{t} \right) + \gamma \mathbb{E}_{\mathbf{x}_{t+1}\sim p} \left[V_{\bar{\theta}} \left(\mathbf{x}_{t+1} \right) \right] \right) \right)^{2} \right],$$
(15)

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:

$$V_{\bar{\theta}}\left(\mathbf{x}_{t}\right) = \mathbb{E}_{\mathbf{a}_{t} \sim \pi}\left[Q_{\bar{\theta}}\left(\mathbf{x}_{t}, \mathbf{a}_{t}\right) - \alpha \log \pi\left(\mathbf{a}_{t} \mid \mathbf{x}_{t}\right)\right].$$
 (16)

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 θ , which helps stabilize training. The soft Bellman residual $J_Q(\theta)$ in Eq. (15) can be optimized with stochastic gradients

$$\hat{\nabla}_{\theta} J_{Q}(\theta) = \nabla_{\theta} Q_{\theta} \left(\mathbf{a}_{t}, \mathbf{x}_{t} \right) \left(Q_{\theta} \left(\mathbf{x}_{t}, \mathbf{a}_{t} \right) - \left(r \left(\mathbf{x}_{t}, \mathbf{a}_{t} \right) + \gamma \left(Q_{\bar{\theta}} \left(\mathbf{x}_{t+1}, \mathbf{a}_{t+1} \right) - \alpha \log \left(\pi_{\phi} \left(\mathbf{a}_{t+1} \mid \mathbf{x}_{t+1} \right) \right) \right) \right) \right).$$
(17)

The policy parameters ϕ can be learned by directly minimizing the expected KL divergence in

$$J_{\pi}(\phi) = \mathbb{E}_{\mathbf{x}_{t} \sim \mathcal{D}} \left[\mathbb{E}_{\mathbf{a}_{t} \sim \pi_{\phi}} \left[\alpha \log \left(\pi_{\phi} \left(\mathbf{a}_{t} \mid \mathbf{x}_{t} \right) \right) - Q_{\theta} \left(\mathbf{x}_{t}, \mathbf{a}_{t} \right) \right] \right].$$
(18)

A neural network transformation is used to parameterize the policy as $\mathbf{a}_t = f_{\phi}(\epsilon_t; \mathbf{x}_t)$, where ϵ_t is an input noise vector sampled from a Gaussian distribution. The objective stated by Eq. (18) can be rewritten as

$$J_{\pi}(\phi) = \mathbb{E}_{\mathbf{x}_{t} \sim \mathcal{D}, \epsilon_{t} \sim \mathcal{N}} \left[\alpha \log \pi_{\phi} \left(f_{\phi} \left(\epsilon_{t}; \mathbf{x}_{t} \right) \mid \mathbf{x}_{t} \right) - Q_{\theta} \left(\mathbf{x}_{t}, f_{\phi} \left(\epsilon_{t}; \mathbf{x}_{t} \right) \right) \right],$$
(19)

where π_{ϕ} is defined implicitly in terms of f_{ϕ} . The gradient of Eq. (19) is approximated with

$$\hat{\nabla}_{\phi} J_{\pi}(\phi) = \nabla_{\phi} \alpha \log \left(\pi_{\phi} \left(\mathbf{a}_{t} \mid \mathbf{x}_{t} \right) \right) + \left(\nabla_{\mathbf{a}_{t}} \alpha \log \left(\pi_{\phi} \left(\mathbf{a}_{t} \mid \mathbf{x}_{t} \right) \right) \\ - \nabla_{\mathbf{a}_{t}} Q \left(\mathbf{x}_{t}, \mathbf{a}_{t} \right) \right) \nabla_{\phi} f_{\phi} \left(\epsilon_{t}; \mathbf{x}_{t} \right),$$
(20)

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. (15) and the policy improvement Eq. (18) converges to the optimal policy. *Proof in* [11].

C. Action Quantization and Correction

The SAC learning algorithm produces the SAC action \mathbf{a}_t at time t that maximizes the policy value $\pi_{\phi} (\mathbf{a}_t | \mathbf{x}_t)$ given the SAC state \mathbf{x}_t . However, to evaluate the reward and update the cache, we need to convert the continuous SAC action \mathbf{a}_t into a discrete action $(c_{A(t)}, \mathbf{b}, \Delta \mathbf{s})$. We achieve this through a simple action quantization method that involves thresholding and integer projection.

Action quantization: Consider an element $\bar{\eta}$ in the SAC action a and its corresponding quantized version η in the system action, where η belongs to the selection set S_{η} . To convert $\bar{\eta}$ to η , we employ uniform thresholding for integer projection, which is given by

$$\eta = \min S_{\eta} + (\bar{\eta} - \min S_{\eta}) \mod \frac{\max S_{\eta} - \min S_{\eta}}{\max S_{\eta} - \min S_{\eta} + 1}$$
(21)

As an example, consider the push action $b_f(t) \in S_{\underline{b}_f} = \{0, 1\}$, we can determine its quantized value $b_f(t) = \overline{b}_f(t) \mod 0.5$ using Eq. (21).

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 λ_Q, λ_{π} , and weight ξ . 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, calculate $r (\mathbf{x}_t, \mathbf{a}_t)$ $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{x}_t, \mathbf{a}_t, r (\mathbf{x}_t, \mathbf{a}_t), \mathbf{x}_{t+1})\}$ end for for each gradient step do $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q (\theta_i)$ for $i \in \{1, 2\}$ $\frac{\phi}{\theta_i} \leftarrow \xi \theta_i + (1 - \xi) \overline{\theta}_i$ for $i \in \{1, 2\}$ end for end for

Action correction: Due to the constraints outlined in Eqs. (2), (3), (8), (9), and (10), the valid action space of the system is very limited and sparsely spanned, with a cardinality of $(M + 1) \times 2^F \times 3^{2F}$. Consequently, even with techniques such as penalty reward, it is difficult for the SAC algorithm to learn which actions are valid in the huge and sparsely spanning action space. Therefore, the post-quantization action $(c_{A(t)}, \mathbf{b}, \Delta \mathbf{s})$ is usually not valid. To overcome this issue, we propose to correct the output action of SAC and make it valid using *Rules 1, 5, and 7* as detailed below. These rules are designed to satisfy the constraints presented in Section II-C. Additionally, we suggest some general *Rules 2, 3, 4, and 6* to accelerate training and refine system action.

- *Rule 1*: when $S_{A(t)}^{O} = 0$, if the $c_{A(t)}$ is smaller than the minimum workable value $\lceil I_{A(t)} w_{A(t)} / (\tau f_D) \rceil$, we will correct it by $c_{A(t)} = \lceil I_{A(t)} w_{A(t)} / (\tau f_D) \rceil$; when $S_{A(t)}^{O} = 0$, $c_{A(t)} = 0$. This is to constrain the total latency.
- Rule 2: when $S_f^I + S_f^O \ge 1$, $b_f = 0$. It indicates no need for proactive pushing if any data of a task is cached.
- *Rule 3*: there is at most one task being proactively transmitted for the task with largest \bar{b}_f and $b_f = 1$. Other tasks are corrected to $b_f = 0$. This is to reduce the unnecessary pushing cost given that the mobile device has one task request at each time slot.
- *Rule* 4: if $b_f = 1$, $\Delta s_f^I = 1$. It indicates that the proactive pushing data has to be cached.
- *Rule 5*: if the cache sum Eq. (10) exceeds the capacity, we drop the input or output cache according to the ascending order of the \bar{s} values until the capacity fits.
- *Rule 6*: if the cache sum Eq. (10) is inferior to capacity, we try to add the reactive input or output cache according to the descending order of $\Delta \bar{s}^{I}_{A(t)}$ and $\Delta \bar{s}^{O}_{A(t)}$ values.
- *Rule* 7: clip the cache action Δs according to the min, max limit in Eq. (8) Eq. (9).



Fig. 2: Training convergence of SAC algorithm for different configuration cases.

D. Reward Design

The reward $r(\mathbf{x}, \mathbf{a})$ of SAC state \mathbf{x} and action \mathbf{a} are designed to be a function of resulting bandwidth and computation cost

$$r(\mathbf{x}, \mathbf{a}) = -\kappa (B(t) + \lambda E(t)) \tag{22}$$

where κ is the normalization coefficient, and is set as 10^{-6} .

The complete SAC learning algorithm is presented in Algorithm 1, where λ_Q, λ_{π} are the step sizes (or learning rate) for stochastic gradient descent, and are chosen to be 1×10^{-4} . ξ is the target smoothing coefficient chosen to be 0.005.

V. IMPLEMENTATION

We generated simulated data for training and testing by randomly generating a Markov chain from the task set \mathcal{F} . The transition probability of a task *i* to one randomly selected task $j \in \mathcal{F} \setminus i$ was set to the maximum transition probability, i.e., $p_{i,j} = p_{\max}$. The probability to other tasks $k \in \mathcal{F} \setminus j$ was given by $p_{i,k} = (1-p_{i,j}) \frac{|p'_{i,k}|}{\sum_{f \in \mathcal{F} \setminus j} |p'_{i,f}|}$, where $p'_{i,k}$ or $p'_{i,f}$ were random samples from a uniform distribution. This designed Markov chain represents the request popularity and transition preference of *F* tasks. We sampled 10⁶ requested tasks using a frame-by-frame approach. In our simulation, we considered M = 8, F = 4, a maximum transition probability of 0.7, $\lambda = 1, I_f = 16000$ bits, $O_f = 30000$ bits, w = 800 cycles/bit, $\tau = 0.02$ seconds, $f_D = 1.7 \times 10^8$ cycles/s, $\mu = 10^{-19}$, and $C = 40 \times 10^3$ bits.

For ease of training and stabilization, both the SAC action \mathbf{a}_t and system state \mathbf{x}_t are normalized to the range of [-1, 1]. The implementation of the system is done with Python and PyTorch. The training and testing were deployed on a PC with TITAN RTX GPU using batch size 256, discount factor $\gamma = 0.99$, automatic entropy temperature α tuning [11], network hidden-layer size 256, one model update per step, one target update per 1000 steps, and replay buffer size of 1×10^7 . We make 10 testing epochs after every 10 training epochs and stop the training and testing when the reward and loss converge.

VI. EVALUATION AND ANALYSIS

A. Baselines

The proposed system is built on the *proactive transmission* and dynamic-computing-frequency reactive service with cache (**PTDFC**). For comparison, we consider the following baselines:

- Most-recently-used proactive transmission and leastrecently-used cache replacement (MRU-LRU): A heuristic algorithm [3] that serves the requested task reactively while proactively caching the *input data* of the most-recently-used task and replacing the least-recentlyused task's cache when the cache is full. The number of computing cores used is fixed at 0.75M.
- Most-frequently-used proactive transmission and leastfrequently-used cache replacement (MFU-LFU): Similar to MRU-LRU, this algorithm replaces the most/least recently used task with the most/least frequently used task.
- Dynamic-computing-frequency reactive service with no cache (DFNC): This algorithm reactively serves the requested task at each time slot t by downloading the input data from the MEC server and computing the output data.
- Dynamic-computing-frequency reactive service with cache (DFC): Similar to DFNC, this algorithm also reactively serves the requested task at each time slot t but can cache the input or output data into limited capacity.

B. Convergence

We show the convergence of the SAC algorithm in Fig. 2 by plotting the reward vs. epochs curves for different system setups. For setups with a smaller action space, such as those with smaller τ and no proactive transmission, the SAC algorithm converges quickly in around 200 epochs. However, more complex setups with larger action spaces, such as those with proactive transmission and more tasks F, typically take 500 epochs or longer to converge.

C. Different Cache Size C

The average transmission cost and computation cost of three SAC-enabled algorithms, DFNC, DFC, PTDFC, and two heuristic algorithms MRU-LRU, MFU-LFU, were compared under different cache sizes C, as shown in Fig. 3. While the cache size change did not affect the DFNC algorithm, the rest of the algorithms showed decreasing transmission costs with increasing C due to the ability to cache more input data locally. In addition, the proposed PTDFC algorithm consistently achieved lower transmission and computation costs than the other algorithms under different cache sizes. For very large cache sizes, such as C = 50000 bits, the performance of PTDFC and DFC was similar.

D. Different Tolerable Service Delays τ

Fig. 4 illustrates the performance of five algorithms at different maximum tolerable service delays. As τ increases, the cost of all algorithms decreases because there is more time available for the transmission and computation process, requiring less bandwidth and lower computing frequency. Among the three SAC-enabled algorithms, the proposed PTDFC algorithm consistently achieves the lowest transmission and computation



Fig. 3: Varying cache size C under default configuration.



Fig. 4: Varying maximum tolerable service latency τ under default configuration.

cost under different τ values. However, as τ gets larger, the transmission cost of all five algorithms begins to converge, and the advantages of PTDFC become less significant. This is because enough time is available for transmission even with the lowest-frequency reactive computing.

VII. CONCLUSION

This paper investigates joint computing, pushing, and caching optimization in a single-user single-server MEC network to reduce the transmission data load and computation cost. A framework based on SAC learning with action quantization and correction techniques is proposed to enable dynamic orchestration of the three activities. Simulation results demonstrate the effectiveness of the proposed framework in reducing both transmission load and computing cost, outperforming baseline algorithms under various parameter settings.

REFERENCES

- Y. Sun, Z. Chen, and H. Liu, "Delay analysis and optimization in cacheenabled multi-cell cooperative networks," in *IEEE Global Communications Conference (GLOBECOM)*, 2016, pp. 1–7.
- [2] 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
- [3] 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.
- [4] W. Chen and H. V. Poor, "Content pushing with request delay information," *IEEE Transactions on Communications*, vol. 65, no. 3, pp. 1146–1161, 2017.
- [5] 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.
- [6] 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.

- [7] 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.
- [8] 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.
- [9] 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. 11365–11373, 2018.
- [10] 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.
- [11] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Offpolicy maximum entropy deep reinforcement learning with a stochastic actor," in *International conference on machine learning*. PMLR, 2018, pp. 1861–1870.