

# Online Reinforcement Learning of X-Haul Content Delivery Mode in Fog Radio Access Networks

Jihwan Moon , *Member, IEEE*, Osvaldo Simeone , *Fellow, IEEE*, Seok-Hwan Park , *Member, IEEE*, and Inkyu Lee , *Fellow, IEEE*

**Abstract**—We consider a Fog Radio Access Network (F-RAN) with a Base Band Unit (BBU) in the cloud and multiple cache-enabled enhanced Remote Radio Heads (eRRHs). The system aims at delivering contents on demand with minimal average latency from a time-varying library of popular contents. Uncached requested files can be transferred from the cloud to the eRRHs by following either backhaul or fronthaul modes. The backhaul mode transfers fractions of the requested files, while the fronthaul mode transmits quantized baseband samples as in Cloud-RAN (C-RAN). The backhaul mode allows the caches of the eRRHs to be updated, which may lower future delivery latencies. In contrast, the fronthaul mode enables cooperative C-RAN transmissions that may reduce the current delivery latency. Taking into account the trade-off between current and future delivery performance, this letter proposes an adaptive selection method between the two delivery modes to minimize the long-term delivery latency. Assuming an unknown and time-varying popularity model, the method is based on model-free Reinforcement Learning (RL). Numerical results confirm the effectiveness of the proposed RL.

**Index Terms**—Caching, F-RAN (Fog RAN), machine learning, reinforcement learning, X-haul.

## I. INTRODUCTION

THE architecture of the recently launched fifth generation (5G) mobile system [1] can leverage cloud processing at Base Band Units (BBUs), as well as edge processing, including edge caching, at enhanced Remote Radio Heads (eRRHs) [2]. In order to enable a flexible functional split in this architecture, referred to as Fog-Radio Access Network (F-RAN) [3], the concept of *X-haul* has been introduced to integrate the traditionally distinct backhaul and fronthaul connectivity modes for the interface between the BBU and the eRRH into a unified framework [4]–[6]. The backhaul mode enables the transfer of data packets from the BBU in the cloud to the eRRHs. In

Manuscript received May 28, 2019; accepted July 5, 2019. Date of publication August 7, 2019; date of current version August 28, 2019. This work was supported by the National Research Foundation through the Ministry of Science, ICT, and Future Planning (MSIP), Korean Government under Grant 2017R1A2B3012316. The work of O. Simeone was supported by the European Research Council (ERC) under the European Union’s Horizon 2020 Research and Innovation Programme under Grant 725731. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Qing Ling. (*Corresponding author: Inkyu Lee.*)

J. Moon and I. Lee are with the School of Electrical Engineering, Korea University, Seoul 02841, South Korea (e-mail: anschino@korea.ac.kr; inkyu@korea.ac.kr).

O. Simeone is with the Department of Informatics, King’s College London, London WC2R 2LS, U.K. (e-mail: osvaldo.simeone@kcl.ac.uk).

S.-H. Park is with the Division of Electronic Engineering, Chonbuk National University, Jeonju 54896, South Korea (e-mail: seokhwan@jbnu.ac.kr).

Digital Object Identifier 10.1109/LSP.2019.2932193

contrast, the fronthaul mode allows the BBU to carry out joint baseband processing and deliver quantized baseband samples to the eRRHs as in Cloud-RAN (C-RAN) [7]–[10].

In this work, we study an adaptive selection of backhaul and fronthaul transfer modes with the aim of optimizing the performance of content delivery. The content delivery in F-RANs has been widely studied in recent years [11]–[17]. Most studies assume offline caching with a static popularity model. Under these assumptions, references [11] and [12] investigated the problem of instantaneous delivery latency minimization and minimum data rate maximization, respectively, while keeping the contents of the caches fixed. In contrast, in [13] and [14], information-theoretic performance bounds were provided on the optimal high Signal-to-Noise-Ratio (SNR) performance by considering also the optimization of uncoded caching strategies. An extension of this work that accounts for time-varying and possibly unknown file popularity with online caching was described in [15]. Under an unknown dynamic popularity model, the works [16] and [17] presented a Reinforcement Learning (RL) based optimization of online caching by assuming a backhaul mode.

In this letter, we investigate for the first time the online minimization of the long-term delivery latency over X-haul links in an F-RAN with time-varying unknown file popularity. We focus on the joint optimization of linear precoding strategies and the choice between fronthaul and backhaul modes. The backhaul mode enables cache updates at the eRRHs, hence potentially reducing future latencies. In contrast, the fronthaul mode allows cooperative C-RAN transmissions which decrease the current delivery latency [11]–[13]. We propose a new model-free RL approach based on a linear value function approximation with properly selected features, and numerical results confirm the effectiveness of the proposed RL scheme.

**Notations:**  $\mathbb{E}[\cdot]$  and  $\Pr(\cdot)$  stand for expectation and probability, respectively.  $|\mathcal{A}|$  represents the cardinality of set  $\mathcal{A}$ , and  $\mathbb{C}^{m \times n}$  denotes an  $m \times n$  complex matrix.  $\mathbb{I}\{c\}$  outputs one if condition  $c$  is true and zero otherwise. For a matrix  $\mathbf{X}$ ,  $|\mathbf{X}|$ ,  $\mathbf{X}^T$ ,  $\mathbf{X}^H$ ,  $\mathbf{X}^{-1}$  and  $\text{tr}(\mathbf{X})$  are defined as determinant, transpose, Hermitian, inverse and trace, respectively.  $\mathbf{I}_m$  means an  $m \times m$  identity matrix while  $\otimes$  equals a Kronecker product operation. Also,  $\text{diag}(\mathbf{X}_1, \dots, \mathbf{X}_N)$  represents block-wise diagonalization of matrices  $\mathbf{X}_1, \dots, \mathbf{X}_N$ . Lastly,  $\mathcal{CN}(\boldsymbol{\mu}, \boldsymbol{\Omega})$  indicates a circularly symmetric complex Gaussian distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Omega}$ .

## II. SYSTEM MODEL

We study the F-RAN system illustrated in Fig. 1, which consists of a BBU in the cloud, connected to  $M$  cache-enabled eRRHs and  $K$  users. Each X-haul link between the BBU and

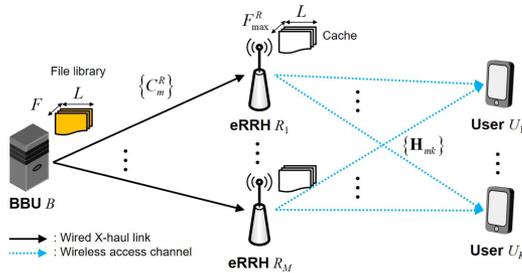


Fig. 1. Illustration of the F-RAN system under study.

the  $m$ -th eRRH has capacity  $C_m^R$  bits per symbols and can be operated in both backhaul and fronthaul modes [5] [6]. The  $k$ -th user and the  $m$ -th eRRH are equipped with  $N_k^U$  and  $N_m^R$  antennas, respectively. We assume a time-slotted operation [17], and the wireless channel matrix  $\mathbf{H}_{mk}$  between the  $m$ -th eRRH and the  $k$ -th user is assumed to be fixed for the given time scale of interest  $T_B$  slots. We also define  $\mathcal{F} \triangleq \{f_1, \dots, f_F\}$  as the library of  $F$   $L$ -bit files, which may be requested by the users. Finally, we denote  $\mathcal{F}^R(t) \subseteq \mathcal{F}$  as the subset of files cached at time slot  $t$  at the eRRHs whose cardinality is bounded by  $F_{\max}^R$  files due to storage capacity constraints. Note that in this letter, we make a simplifying assumption that all the eRRHs store the same files in their respective caches. Generalization of the framework is possible but at the cost of a more cumbersome notation.

#### A. Request Model and Online Caching

In each time slot  $t$ , a subset  $\mathcal{F}_{\text{pop}}(t) \subseteq \mathcal{F}$  of files is popular in the sense that all users request files from  $\mathcal{F}_{\text{pop}}(t)$ . Specifically, the  $k$ -th user requests a uniformly selected file  $f_k^U(t)$  from subset  $\mathcal{F}_{\text{pop}}(t)$  without replacement [15]. The assumption of no replacement ensures that all requested files are distinct, yielding a worst-case performance analysis [13]. We assume that the popularity  $\mathcal{F}_{\text{pop}}(t)$  varies as a Markov process as in [16], [18]–[20]. This is a standard assumption which provides a first-order approximation of the evolution of the content popularity [21][22]. Let  $\mathcal{K}_{\text{req,C}}(t)$  and  $\mathcal{K}_{\text{req,NC}}(t)$  denote the indices of the users whose requested files  $\mathcal{F}_{\text{req,C}}(t) \triangleq \{f_k^U(t)\}_{k \in \mathcal{K}_{\text{req,C}}(t)}$  are cached and the indices of users whose requested files  $\mathcal{F}_{\text{req,NC}}(t) \triangleq \{f_k^U(t)\}_{k \in \mathcal{K}_{\text{req,NC}}(t)}$  are not cached at time  $t$ , respectively. In case the backhaul mode is selected at time slot  $t$ , the requested but uncached files in  $\mathcal{F}_{\text{req,NC}}(t)$  are sent on all the X-haul links and cached. In order to make space for a new file, a previously cached file is evicted by following the standard Least Recently Used (LRU) rule [23].

#### B. Delivery Operation

At each slot  $t$ , the X-haul link is used in either fronthaul or backhaul mode for  $\Delta^R(t, a(t))$  symbols, where  $a(t) = 0$  and 1 indicate the selection of fronthaul and backhaul modes, respectively. Subsequently, the eRRHs deliver the requested files in set  $\mathcal{F}_{\text{req}}(t) \triangleq \mathcal{F}_{\text{req,C}}(t) \cup \mathcal{F}_{\text{req,NC}}(t)$  over the wireless channel for  $\Delta^U(t, a(t))$  symbols, based on the signals received on the X-haul links and on the cached contents. This results in a total latency of  $\Delta(t, a(t)) = \Delta^R(t, a(t)) + \Delta^U(t, a(t))$  symbols for time slot  $t$ . Note that the eRRHs' caches are updated according to the caching mechanism described in Section II-A only if the backhaul mode is selected as  $a(t) = 1$ .

#### C. Problem Formulation

The delivery time  $\Delta(t, a(t))$  at slot  $t$  depends on the state of the system  $s(t) = \{\mathcal{F}_{\text{pop}}(t), \mathcal{F}^R(t), \mathcal{F}_{\text{req}}(t)\}$ , which includes the set of popular files, cached files and requested files, respectively. Given the Markovity of the process  $\mathcal{F}_{\text{pop}}(t)$ , the state  $s(t)$  evolves as a controlled Markov process.  $s(t)$  is partially observable since the set  $\mathcal{F}_{\text{pop}}(t)$  is unknown, and it is only observed indirectly via the file set  $\mathcal{F}_{\text{req}}(t)$ . In particular, at time  $t$ , only the history of observations  $\mathbf{o}(1:t) \triangleq \{\mathbf{o}(1), \dots, \mathbf{o}(t)\}$  with  $\mathbf{o}(t) = \{\mathcal{F}_{\text{req}}(t), \mathcal{F}^R(t)\}$  is available to the system. Thus, a general policy can map the observations  $\mathbf{o}(1:t)$  to the selected action  $a(t)$  through a conditional distribution  $\pi(a(t)|\mathbf{o}(1:t))$ .

In this work, we aim at minimizing the average long-term delivery latency of the proposed F-RAN system over the selection of policy  $\pi(a(t)|\mathbf{o}(1:t))$ . Given a forgetting factor  $\gamma \leq 1$ , the problem can be formulated as

$$(P): \min_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=1}^{\infty} \gamma^t \Delta(t, a(t)) \right] \text{ s.t. } a(t) \in \{0, 1\}, \forall t, \quad (1)$$

where calculation of the total latency  $\Delta(t, a(t))$  will be reviewed in Section III. The expectation in (P) is over the state distribution, which depends on the policy.

### III. MINIMUM INSTANTANEOUS LATENCY

In this section, we discuss how to evaluate the delivery latency  $\Delta(t, a(t))$  in problem (P). We emphasize that  $\Delta(t, a(t))$  for  $a(t) = 0$  and 1 is assumed known when solving problem (P) at each time slot  $t$ , and is derived as defined in this section. Following [11], we omit the time index  $t$  for simplicity.

#### A. Backhaul Mode

In the backhaul mode ( $a = 1$ ), the BBU first fetches the requested but uncached files  $\mathcal{F}_{\text{req,NC}}$  and transmits them to the eRRHs. The backhaul transmission to the  $m$ -th eRRH takes  $\Delta_m^R = |\mathcal{F}_{\text{req,NC}}|L/C_m^R$  symbols, and the total backhaul latency is  $\Delta^R = \max_m \Delta_m^R$ , since all the eRRHs need to receive the files in  $\mathcal{F}_{\text{req,NC}}$ . As a result, all the requested files in  $\mathcal{F}_{\text{req}}$  are available at the eRRHs and cooperative transmission across all eRRHs is feasible. Each file  $f_k^U \in \mathcal{F}_{\text{req}}$  for the  $k$ -th user is encoded by each eRRH as the signal  $\mathbf{s}_k \in \mathbb{C}^{n_k \times 1} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_{n_k})$ , where  $n_k \leq N_k^U$  denotes the number of data streams allocated to the  $k$ -th user, which is assumed to be a fixed parameter. The transmit signal from the  $m$ -th eRRH is then given as  $\mathbf{x}_m = \sum_{k \in \mathcal{K}_{\text{req}}} \mathbf{G}_{mk} \mathbf{s}_k$  where  $\mathcal{K}_{\text{req}} \triangleq \mathcal{K}_{\text{req,C}} \cup \mathcal{K}_{\text{req,NC}}$ , and  $\mathbf{G}_{mk} \in \mathbb{C}^{N_m^R \times n_k}$  is the precoding matrix for  $\mathbf{s}_k$  at the  $m$ -th eRRH. Accordingly, the achievable rate for the  $k$ -th user on the wireless channel can be written as [11]

$$R_{\text{back},k}^U(\{\mathbf{G}_k\}) = \log_2 \left| \mathbf{I}_{N_k^U} + \Phi_{\text{back},k}^U \right| \text{ [bits/symbol]}, \quad (2)$$

where we have  $\Phi_{\text{back},k}^U \triangleq (\sum_{\ell \in \mathcal{K}_{\text{req}} \setminus k} \mathbf{H}_k \mathbf{G}_{\ell} \mathbf{G}_{\ell}^H \mathbf{H}_k^H + \sigma_k^2 \mathbf{I}_{N_k^U})^{-1} \mathbf{H}_k \mathbf{G}_k \mathbf{G}_k^H \mathbf{H}_k^H$  with  $\mathbf{H}_k \triangleq [\mathbf{H}_{1k} \cdots \mathbf{H}_{Mk}]$  and  $\mathbf{G}_k \triangleq [\mathbf{G}_{1k}^T \cdots \mathbf{G}_{Mk}^T]^T$ , and  $\sigma_k^2$  represents the additive white Gaussian noise variance at the  $k$ -th user.

The latency  $\Delta_k^U$  for delivering file  $f_k^U$  for the  $k$ -th user is obtained as  $\Delta_k^U = L/R_{\text{back},k}^U(\{\mathbf{G}_k\})$ , and the overall wireless channel latency equals  $\Delta^U = \max_k \Delta_k^U$ , since every requesting user needs to receive the requested file. The minimum instantaneous latency  $\Delta$  for  $a = 1$  can hence be found as a solution of

the problem

$$(P1): \min_{\Delta^U, \{\mathbf{G}_k\}} \Delta^R + \Delta^U \quad (3a)$$

$$\text{s.t. } \Delta^U \geq L/R_{\text{back},k}^U(\{\mathbf{G}_k\}), \forall k \in \mathcal{K}_{\text{req}}, \quad (3b)$$

$$\text{tr} \left( \sum_{k \in \mathcal{K}_{\text{req}}} \mathbf{E}_m \mathbf{G}_k \mathbf{G}_k^H \mathbf{E}_m^H \right) \leq P_m^R, m = 1, \dots, M, \quad (3c)$$

where  $P_m^R$  denotes the maximum transmit power of the  $m$ -th eRRH, and we define  $\mathbf{E}_m \triangleq [\mathbf{0} \cdots \mathbf{I}_{N_m^R} \cdots \mathbf{0}]$  in which an identity matrix  $\mathbf{I}_{N_m^R}$  spans columns from  $\sum_{\ell=1}^{m-1} N_\ell^R + 1$  to  $\sum_{\ell=1}^m N_\ell^R$ . Although problem (P1) is jointly non-convex, a stationary point can be attained by leveraging Successive Convex Approximation (SCA) as detailed in [11].

### B. Fronthaul Mode

Under the fronthaul mode, any requested but uncached file  $f_k^U \in \mathcal{F}_{\text{req,NC}}$  for the  $k$ -th user is jointly encoded and precoded at the BBU. The resulting signal dedicated for the  $m$ -th eRRH is written as  $\hat{\mathbf{x}}_m = \sum_{k \in \mathcal{K}_{\text{req,NC}}} \mathbf{W}_{mk} \mathbf{s}_k$ , where  $\mathbf{s}_k \in \mathbb{C}^{n_k \times 1} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_{n_k})$  encodes file  $f_k^U$ , and  $\mathbf{W}_{mk} \in \mathbb{C}^{N_m^R \times n_k}$  represents the corresponding precoding matrix for the  $m$ -th eRRH. The BBU then performs compression on  $\hat{\mathbf{x}}_m$  prior to transferring to the eRRHs. As a result, the decompressed signal at the  $m$ -th eRRH can be written by  $\tilde{\mathbf{x}}_m = \hat{\mathbf{x}}_m + \mathbf{q}_m$  with quantization noise  $\mathbf{q}_m \in \mathbb{C}^{N_m^R \times 1} \in \mathcal{CN}(\mathbf{0}, \mathbf{\Omega}_m)$  for a given covariance matrix  $\mathbf{\Omega}_m$  [11] [12].

The rest of the requested cached files  $\mathcal{F}_{\text{req,C}}$  are locally precoded with  $\{\mathbf{G}_{mk}\}$  at the eRRHs in the same manner as in the backhaul mode. The final transmit signal at the  $m$ -th eRRH is then given as  $\mathbf{x}_m = \sum_{k \in \mathcal{K}_{\text{req,C}}} \mathbf{G}_{mk} \mathbf{s}_k + \tilde{\mathbf{x}}_m$ , and the achievable rate for the  $k$ -th user can be obtained as [11]

$$R_{\text{front},k}^U(\{\tilde{\mathbf{G}}_k\}, \mathbf{\Omega}_R) = \log_2 \left| \mathbf{I}_{N_k^U} + \mathbf{\Phi}_{\text{front},k}^U \right| \text{ [bits/symbol]}, \quad (4)$$

where we have  $\mathbf{\Phi}_{\text{front},k}^U \triangleq (\sum_{\ell \in \mathcal{K}_{\text{req}} \setminus k} \mathbf{H}_k \tilde{\mathbf{G}}_\ell \tilde{\mathbf{G}}_\ell^H \mathbf{H}_k^H + \mathbf{H}_k \mathbf{\Omega}_R \mathbf{H}_k^H + \sigma_k^2 \mathbf{I}_{N_k^U})^{-1} \mathbf{H}_k \tilde{\mathbf{G}}_k \tilde{\mathbf{G}}_k^H \mathbf{H}_k^H$ ,  $\mathbf{\Omega}_R \triangleq \text{diag}(\mathbf{\Omega}_1, \dots, \mathbf{\Omega}_M)$ ,  $\tilde{\mathbf{G}}_k \triangleq [\tilde{\mathbf{G}}_{1k}^T \cdots \tilde{\mathbf{G}}_{Mk}^T]^T$  with  $\tilde{\mathbf{G}}_{mk} \triangleq b_k^U \mathbf{G}_{mk} + (1 - b_k^U) \mathbf{W}_{mk}$ , and  $b_k^U = 1$  if  $f_k^U \in \mathcal{K}_{\text{req,C}}$  and  $b_k^U = 0$  otherwise for the  $k$ -th user.

The wireless channel latency  $\Delta^U$  is defined in the same way as in the backhaul mode. For the fronthaul latency, by the rate-distortion theory, sending quantized signals to the  $m$ -th eRRH consumes

$$g_m(\{\tilde{\mathbf{G}}_k\}, \mathbf{\Omega}_R) = \log_2 \left| \mathbf{I}_{N_m^R} + \mathbf{\Phi}_m^R \right| \text{ [bits/symbol]}, \quad (5)$$

with  $\mathbf{\Phi}_m^R \triangleq (\mathbf{E}_m \mathbf{\Omega}_R \mathbf{E}_m^H)^{-1} \sum_{k \in \mathcal{K}_{\text{req,NC}}} \mathbf{E}_m \tilde{\mathbf{G}}_k \tilde{\mathbf{G}}_k^H \mathbf{E}_m^H$  [11]. Compressing  $\Delta^U$  symbols produces  $\Delta^U g_m(\{\tilde{\mathbf{G}}_k\}, \mathbf{\Omega}_R)$  bits, which need to be transferred from the BBU to the  $m$ -th eRRH. Therefore, the fronthaul latency is given by  $\Delta^R = \max_m \Delta_m^R$  where  $\Delta_m^R = \Delta^U g_m(\{\tilde{\mathbf{G}}_k\}, \mathbf{\Omega}_R) / C_m^R$ , and the minimum instantaneous latency  $\Delta$  for  $a = 0$  is calculated as a solution of the problem

$$(P2): \min_{\Delta^R, \Delta^U, \{\tilde{\mathbf{G}}_k\}, \mathbf{\Omega}_R} \Delta^R + \Delta^U \quad (6a)$$

$$\text{s.t. } \Delta^R \geq \Delta^U g_m(\{\tilde{\mathbf{G}}_k\}, \mathbf{\Omega}_R) / C_m^R, m = 1, \dots, M, \quad (6b)$$

$$\Delta^U \geq L/R_{\text{front},k}^U(\{\tilde{\mathbf{G}}_k\}, \mathbf{\Omega}_R), \forall k \in \mathcal{K}_{\text{req}}, \quad (6c)$$

$$\text{tr} \left( \sum_{k \in \mathcal{K}_{\text{req}}} \mathbf{E}_m \tilde{\mathbf{G}}_k \tilde{\mathbf{G}}_k^H \mathbf{E}_m^H + \mathbf{E}_m \mathbf{\Omega}_R \mathbf{E}_m^H \right) \leq P_m^R, \quad (6d)$$

$$m = 1, \dots, M,$$

which can be tackled via the SCA approach detailed in [11]. The total worst-case order of complexity for the SCA method can be expressed as  $\mathcal{O}(N_{\text{SCA}} \sqrt{N_{\text{const}}} \log(N_{\text{const}}/\epsilon))$  where  $\epsilon$ ,  $N_{\text{SCA}}$  and  $N_{\text{const}}$  indicate the desired error tolerance, the maximum number of the SCA iterations and the number of constraints, respectively [24]. Here,  $N_{\text{const}}$  equals  $|\mathcal{K}_{\text{req}}| + M$  in (P1) and  $|\mathcal{K}_{\text{req}}| + 2M$  in (P2).

## IV. RL-BASED X-HAUL ONLINE OPTIMIZATION

In this section, we solve problem (P) by proposing an online on-policy RL-based optimization strategy [25].

### A. Problem (P) as a Partially Observable Decision Process

As discussed in Section II, problem (P) is a Partially Observable Markov Decision Process (POMDP) with the action space  $\{0, 1\}$  and the instantaneous reward given by the negative latency  $r(t+1) = -\Delta(t, a(t))$ . In order to reduce the complexity of the policy, we optimize here over memoryless policies that select an action  $a(t)$  based only on the latest observation  $\mathbf{o}(t)$  at time slot  $t$  [26] [27] as well as a summary of the previous observations  $\mathbf{o}(1:t)$  given by the set  $\{\tau_{\text{req},f}(t)\}_{f \in \mathcal{F}^R(t)}$  where  $\tau_{\text{req},f}(t)$  is the most recent time slot at which cached file  $f$  was requested at time slot  $t$ .

### B. SARSA With Linear Value Function Approximation

To optimize over memoryless policies, we adopt the online on-policy value-based strategy State-Action-Reward-State-Action (SARSA) with a carefully designed linear approximation [25]. The SARSA updates an action-value function, or Q-function,  $q(o, a)$  that estimates the expected return  $\mathbb{E}[G(t)|\mathbf{o} = o, a = a]$  with  $G(t) \triangleq \sum_{\tau=0}^{\infty} \gamma^\tau r(t+\tau+1)$ . Since the total size of the observation space in (P) grows exponentially with  $F$ , we propose a linear value function approximation  $\hat{q}(o, a, \mathbf{w}) \triangleq \mathbf{w}^T \mathbf{x}(o, a)$ , where  $\mathbf{w}$  is a parameter vector to be learned, and  $\mathbf{x}(o, a)$  denotes a feature vector representing the observation-action pair  $(o, a)$  [25].

In order to determine a suitable feature vector, we first note that vector  $\mathbf{x}(o, a)$  should contain sufficient information to quantify the value of caching for currently cached and requested files. Frequently requested files typically yield lower future latencies when cached, but an optimal choice should account not only for their popularity but also for their remaining *life time*, which is a duration that a file remains popular (see Sec. II of [28] for further discussion).

Based on these considerations, we introduce a variable  $\phi_\ell(t)$  for every file  $f_\ell \in \mathcal{F}$  as a function of the current observation  $\mathbf{o}(t)$  at time slot  $t$ . We set it as  $\phi_\ell(t) = 1$  if  $f_\ell \in \mathcal{F}_{\text{req,NC}}(t)$ ,  $\phi_\ell(t) = 2$  if  $f_\ell \in \mathcal{F}^R(t)$  and  $\phi_\ell(t) = 0$  otherwise. Furthermore, we also include a variable  $\theta(t) \triangleq t - \max_{f \in \mathcal{F}^R(t)} \tau_{\text{req},f}$  that measures the ‘‘age’’ of the currently cached files, that is, the maximum time elapsed since the last request of the cached files. We can quantize this variable by  $N_\Theta$  ranges  $\Theta_1, \dots, \Theta_{N_\Theta} \subseteq \mathbb{R}^+$  with  $\Theta_i \cap \Theta_j = \emptyset$  for all  $i \neq j$  and  $\bigcup \Theta_i = \mathbb{R}^+$ . If the caches are

**Algorithm 1:** Proposed RL-Based Solution for Problem (P).

---

Initialize the total number of episodes  $N_{\text{epi}}$ , weight vector  $\mathbf{w} = \mathbf{0}$ , eligibility trace  $\mathbf{E} = \mathbf{0}$ , and  $\gamma, \lambda \in (0, 1]$

For  $n_{\text{epi}} = 1 : N_{\text{epi}}$

    Randomly initialize cached contents  $\mathcal{F}^R(0)$  and generate  $\{\mathbf{H}_{mk}\}$

    For  $t = 1 : T_B$

        Collect observation  $\mathbf{o}(t) = \{\mathcal{F}_{\text{req}}(t), \mathcal{F}^R(t), \{\tau_{\text{req},f}(t)\}_{f \in \mathcal{F}^R(t)}\}$

        Choose the delivery mode greedily with probability  $1 - 1/n_{\text{epi}}$  as  $\mathbf{a}(t) = \arg \max_{a'} \mathbf{w}^T \mathbf{x}(\mathbf{o}(t), a')$ , and uniformly with probability  $1/n_{\text{epi}}$

        If  $\mathbf{a}(t) = 1$ , update  $\mathcal{F}_{\text{cache},R}(t)$  according to LRU

        Set  $r(t+1) = -\Delta(t, \mathbf{a}(t))$

        Update  $\mathbf{E} \leftarrow \gamma \lambda \mathbf{E} + \mathbf{x}(\mathbf{o}, a)$

        Update  $\mathbf{w} \leftarrow \mathbf{w} + \beta \delta(t, \mathbf{w}) \mathbf{E}$  with  $\beta = 1/n_{\text{epi}}$

    End

End

---

up to date, the quantity  $t - \tau_{\text{req},f}$  is small for all  $f \in \mathcal{F}^R(t)$ , and hence  $\theta(t)$  is also small. Otherwise, if there exists any file  $f \in \mathcal{F}^R(t)$  with large  $t - \tau_{\text{req},f}$ , a refresh of the caches may be required.

Using the variables introduced above, we define the feature vector  $\mathbf{x}(\mathbf{o}(t), \mathbf{a}(t))$  as

$$\mathbf{x}(\mathbf{o}(t), \mathbf{a}(t)) = [\phi_1^T(t) \cdots \phi_F^T(t) \boldsymbol{\theta}^T(t)]^T \otimes \mathbf{a}(t), \quad (7)$$

where we have used the one-hot encoded vectors  $\phi_\ell(t) \triangleq [\mathbb{I}\{\phi_\ell(t) = 1\} \mathbb{I}\{\phi_\ell(t) = 2\} \mathbb{I}\{\phi_\ell(t) = 0\}]^T$ ,  $\boldsymbol{\theta}(t) \triangleq [\mathbb{I}\{\theta(t) \in \Theta_1\} \cdots \mathbb{I}\{\theta(t) \in \Theta_{N_\Theta}\}]^T$  and  $\mathbf{a}(t) \triangleq [\mathbb{I}\{\mathbf{a}(t) = 0\} \mathbb{I}\{\mathbf{a}(t) = 1\}]^T$ . The feature vector  $\mathbf{x}(\mathbf{o}(t), \mathbf{a}(t))$  in (7) has dimension  $2(N_\Theta + 3F)$ , which increases linearly in  $F$  and is hence significantly smaller than the size of the conventional look-up table-based SARSA. The effectiveness of the proposed feature vector  $\mathbf{x}(\mathbf{o}(t), \mathbf{a}(t))$  will be verified in Section V.

The overall proposed procedure for solving (P) is summarized in Algorithm 1 where  $\delta(t, \mathbf{w}) \triangleq r(t+1) + \gamma \hat{q}(\mathbf{o}(t+1), \mathbf{a}(t+1), \mathbf{w}) - \hat{q}(\mathbf{o}, \mathbf{a}, \mathbf{w})$  denotes the temporal difference error, and  $\mathbf{E}$  indicates the eligibility trace. Here, an  $\epsilon$ -greedy exploration strategy with decreasing  $\epsilon$  is adopted. Note that  $\mathbf{E}$  is used to assign credit for the current reward to the most frequently visited states and selected actions, so as to enable online learning (see [25] for details).

## V. NUMERICAL RESULTS

In this section, the performance of the proposed RL-based algorithm is evaluated via numerical examples. We adopt the channel model  $\mathbf{H}_{mk} = \sqrt{\rho_{mk}} \hat{\mathbf{H}}_{mk}$ , where  $\rho_{mk} \triangleq \rho_0 (\frac{d_{mk}}{d_0})^{-\eta}$  equals the distance-dependent path loss between eRRH  $R_m$  and user  $U_k$ ,  $\rho_0$  indicates the path loss at reference distance  $d_0$ ,  $\eta$  is the path loss exponent, and  $d_{mk}$  represents the distance between the  $m$ -th eRRH and the  $k$ -th user. Each element of  $\hat{\mathbf{H}}_{mk}$  follows an independent complex Gaussian distribution with zero mean and unit variance. The eRRHs and the users are circularly placed from the BBU at the center with uniformly distributed angles and distance  $d_{BR} = 200$  m and  $d_{BU} = 400$  m, respectively. The bandwidth is 20 MHz and the thermal noise is  $-170$  dBm/Hz. We set  $K = 10$ ,  $M = 3$ ,  $\rho_0 = 10^{-3}$ ,  $d_0 = 1$  m,  $\eta = 3.5$ ,  $T_B = 100$  time slots,  $F_{\text{max}}^R = 4$  files,  $P_m^R = 30$  dBm,  $N_m^R = N_k^U = 1$  and  $C_m^R = 0.1$  bits per symbol. For RL, we use the hyperparameters  $\gamma = 1$ ,  $\lambda = 0.5$ ,

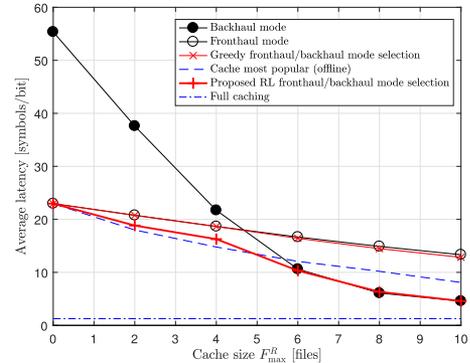


Fig. 2. Average latency with respect to the maximum cache size  $F_{\text{max}}^R$ .

and  $\Theta_\ell = [2(\ell - 1), \min(2(\ell - 1) + 1, \theta_{\text{max}})]$  with  $N_\Theta = 11$  where  $\theta_{\text{max}} = 20$  limits the maximum value of  $\theta(t)$ .

Reference [28] demonstrated that the popularity of files often exhibits temporal locality in the sense that the content is frequently requested in a bursty fashion for a certain life time. Motivated by these findings, we model the evolution of the subset  $\mathcal{F}_{\text{pop}}(t)$  of popular files such that a currently unpopular file  $f$  has a probability of  $P_{\text{pop},f}$  to become popular, and file  $f$  remains popular for  $T_{\text{life},f}$  time slots. We assume Zipf's distribution [29] for  $P_{\text{pop},f} = \ell^{-\xi} / \sum_{\nu=1}^F \nu^{-\xi}$  with  $\xi = 1$ . The proposed RL scheme is compared with a greedy fronthaul/backhaul mode selection that minimizes the current delivery latency at each time slot as well as with an offline scheme that keeps the  $F_{\text{max}}^R$  most popular files with the largest  $P_{\text{pop},f}$  under the idealized assumption that this is known in prior.

Fig. 2 compares the average long-term latency performance as a function of the eRRHs' cache size  $F_{\text{max}}^R$  for  $P_m^R = 30$  dBm,  $T_{\text{life},f} = 10$  and  $F = 20$ . We also limit the maximum number of the SCA iterations for solving (P1) and (P2) as  $N_{\text{SCA}} = 10$ . Note that the convergence to a stationary point for SCA does not affect the convergence of SARSA since we treat the negative reward function  $-\Delta(t, \mathbf{a}(t))$  as fixed. With  $F_{\text{max}}^R \leq 4$ , the fronthaul mode is seen to yield a lower latency than the backhaul mode given the limited advantage of caching in this regime. The opposite is true when the eRRHs have larger caches, such as  $F_{\text{max}}^R > 4$ , in which the backhaul mode outperforms the fronthaul mode. In agreement with the results in [11]–[13] and [15], the greedy scheme almost always selects the fronthaul mode and is hence strongly suboptimal for large enough  $F_{\text{max}}^R$ . The proposed RL method exhibits the lowest latency among all schemes that do not assume the knowledge of the popularity probability. It can be checked that the gain is not obtained by statically selecting the best mode at each time instant, but rather by carrying out an optimized dynamic selection. It is also observed that in a large  $F_{\text{max}}^R$  regime, the proposed strategy can outperform the static offline scheme which assumes popularity dynamics to be known in advance.

## VI. CONCLUSION

In this letter, we have demonstrated the advantage of adaptively selecting between the backhaul and fronthaul transfer modes as a function of the current cache contents and the history of past requests in an F-RAN system. The proposed RL-based strategy has been shown via numerical results to outperform baseline schemes, confirming the potential advantages of an X-haul implementation over static fronthaul or backhaul deployments.

## REFERENCES

- [1] G. Wang, G. Feng, S. Qin, and R. Wen, "Efficient traffic engineering for 5G core and backhaul networks," *J. Commun. Netw.*, vol. 19, no. 1, pp. 80–92, Feb. 2017.
- [2] Y.-J. Ku *et al.*, "5G radio access network design with the fog paradigm: Confluence of communications and computing," *IEEE Commun. Mag.*, vol. 55, no. 4, pp. 46–52, Apr. 2017.
- [3] Y.-Y. Shih, W.-H. Chung, A.-C. Pang, T.-C. Chiu, and H.-Y. Wei, "Enabling low-latency applications in fog-radio access networks," *IEEE Netw.*, vol. 31, no. 1, pp. 52–58, Jan. 2017.
- [4] A. D. L. Oliva *et al.*, "Xhaul: Toward an integrated fronthaul/backhaul architecture in 5G networks," *IEEE Wireless Commun.*, vol. 22, no. 5, pp. 32–40, Oct. 2015.
- [5] T. Pfeiffer, "Next generation mobile fronthaul and midhaul architectures," *J. Opt. Commun. Netw.*, vol. 7, pp. 38–45, Nov. 2015.
- [6] N. J. Gomes, P. Chanclou, P. Turnbull, A. Magee, and V. Jungnickel, "Fronthaul evolution: From CPRI to Ethernet," *Opt. Fiber Technol.*, vol. 26, pp. 50–58, Dec. 2015.
- [7] H. Ren *et al.*, "Low-latency C-RAN: an next-generation wireless approach," *IEEE Veh. Technol. Mag.*, vol. 13, no. 2, pp. 48–56, Jun. 2018.
- [8] J. Kim, H. Lee, S.-H. Park, and I. Lee, "Minimum rate maximization for wireless powered cloud radio access networks," *IEEE Trans. Veh. Technol.*, vol. 68, no. 1, pp. 1045–1049, Jan. 2019.
- [9] J. Kim, S.-H. Park, O. Simeone, I. Lee, and S. S. (Shitz), "Joint design of fronthauling and hybrid beamforming for downlink C-RAN systems," *IEEE Trans. Commun.*, vol. 67, no. 6, pp. 4423–4434, Jun. 2019.
- [10] Y. Jeon, S.-H. Park, C. Song, J. Moon, S. Maeng, and I. Lee, "Joint designs of fronthaul compression and precoding for full-duplex cloud radio access networks," *IEEE Wireless Commun. Lett.*, vol. 5 no. 6, pp. 632–635, Dec. 2016.
- [11] S.-H. Park, O. Simeone, and S. Shamai, "Joint cloud and edge processing for latency minimization in fog radio access networks," in *Proc. IEEE Int. Workshop Signal Adv. Wireless Commun.*, Jul. 2016, pp. 1–5.
- [12] S.-H. Park, O. Simeone, and S. Shamai, "Joint optimization of cloud and edge processing for fog radio access networks," *IEEE Trans. Wireless Commun.*, vol. 15, no. 11, pp. 7621–7632, Nov. 2016.
- [13] A. Sengupta, R. Tandon, and O. Simeone, "Fog-aided wireless networks for content delivery: Fundamental latency Tradeoffs," *IEEE Trans. Inf. Theory*, vol. 63, no. 10, pp. 6650–6678, Oct. 2017.
- [14] J. Zhang and O. Simeone, "Fundamental limits of cloud and cache-aided interference management with multi-antenna edge nodes," *IEEE Trans. Inf. Theory*, vol. 65, no. 8, pp. 5197–5214, Aug. 2019.
- [15] S. M. Azimi, O. Simeone, A. Sengupta, and R. Tandon, "Online edge caching and wireless delivery in fog-aided networks with dynamic content popularity," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 6, pp. 1189–1202, Jun. 2018.
- [16] A. Sadeghi, F. Sheikholeslami, and G. B. Giannakis, "Optimal and scalable caching for 5G using reinforcement learning of space-time popularities," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 1, pp. 180–190, Feb. 2018.
- [17] S. O. Somuyiwa, A. Gyrgy, and D. Gunduz, "A reinforcement-learning approach to proactive caching in wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 6, pp. 1331–1344, Jun. 2018.
- [18] R. Pedarsani, M. A. Maddah-Ali, and U. Niesen, "Online coded caching," *IEEE/ACM Trans. Netw.*, vol. 24, no. 2, pp. 836–845, Apr. 2016.
- [19] N. Garg, M. Sellathurai, and T. Ratnarajah, "Content placement learning for success probability maximization in wireless edge caching networks," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, May 2019, pp. 3092–3096.
- [20] A. Sadeghi, G. Wang, and G. B. Giannakis, "Adaptive caching via deep reinforcement learning," [Online]. Available: <https://arxiv.org/abs/1902.10301>.
- [21] M. Zorzi, R. R. Raoz, and L. B. Milstein, "On the accuracy of a first-order Markov model for data transmission on fading channels," in *Proc. IEEE 4th Int. Conf. Universal Pers. Commun.*, Nov. 1995, pp. 211–215.
- [22] G. Carofiglio, M. Gallo, L. Muscariello, and D. Perino, "Modeling data transfer in content-centric networking," in *Proc. 23rd Int. Teletraffic Congr.*, Sep. 2011, pp. 111–118.
- [23] T. Johnson and D. Shasha, "2Q: A low overhead high performance buffer management replacement algorithm," in *Proc. 20th Int. Conf. Very Large Data Bases*, Sep. 1994, pp. 439–450.
- [24] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, Mar. 2004.
- [25] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. 2nd ed. Cambridge, MA, USA: MIT Press, Oct. 2018.
- [26] M. L. Littman, "Memoryless policies: Theoretical limitations and practical results," in *Proc. Int. Conf. Simul. Adaptive Behav.*, Aug. 1994, pp. 238–245.
- [27] Y. Li, B. Yin, and H. Xi, "Finding optimal memoryless policies of POMDPs under the expected average reward criterion," *Eur. J. Oper. Res.*, vol. 211, pp. 556–567, Jun. 2011.
- [28] S. Traverso, M. Ahmed, M. Garetto, P. Giaccone, E. Leonardi, and S. Niccolini, "Temporal locality in today's content caching: Why it matters and how to model it," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 43, pp. 5–12, Oct. 2013.
- [29] D. M. W. Powers, "Applications and explanations of Zipf's law," in *Proc. Joint Conf. New Methods Lang. Process. Comput. Natural Lang. Learn.*, Jan. 1998, pp. 151–160.