# Edge Computing Assisted Adaptive Streaming Scheme for Mobile Networks

**MINSU KIM** AND **KWANGSUE CHUNG**, (Senior Member, IEEE)

Department of Electronics and Communications Engineering, Kwangwoon University, Seoul 01897, South Korea

Corresponding author: Kwangsue Chung (kchung@kw.ac.kr)

**ABSTRACT** Various video streaming platforms are responsible for continuously generating a majority of Internet traffic. The promising solution for the smooth video streaming services is to use the HTTP adaptive streaming (HAS). To guarantee a high quality of experience (QoE), the clients using the HAS dedicate the bitrate adaptation for the upcoming video segments. However, the client-driven approach causes degradation in QoE, inefficient resource utilization, and unfair bitrate allocation when multiple clients stream a video over the same access network. Edge computing-assisted adaptive streaming gives an opportunity to jointly optimize the QoE of clients, resource utilization, and fairness among clients by shifting the adaptation intelligence from the clients to the edge cloud. In this paper, we first present an adaptive streaming framework taking advantage of the capabilities of multi-access edge computing (MEC). Next, we design an optimization model that jointly considers the main influencing factors in QoE and fairness among clients. The proposed scheme formulates the joint optimization problem with the various constraints by considering the correlation between QoE and fairness. To efficiently solve the joint optimization problem, we propose a greedy-based bitrate allocation algorithm for multiple clients. The results from the performance evaluation show that the proposed scheme can improve the QoE of clients and resource utilization compared with the existing schemes and minimize the loss in fairness.

**INDEX TERMS** HTTP adaptive streaming, quality of experience, mobile networks, multi-access edge computing, resource utilization, fairness.

## I. INTRODUCTION

According to the Cisco Visual Networking Index, most of the Internet traffic is generated through videos, which is expected to increase up to 82% of the total traffic by 2022 [1]. Video streaming service platforms such as Netflix and YouTube generate the majority of the video traffic. With the increase in video traffic, providing smooth video streaming services to clients has become an important challenge. Video streaming services require a solution that can adjust the video bitrate delivered to the clients for smooth video playback. The most promising solution for this issue is to use HTTP adaptive streaming (HAS) [2]. The HAS server encodes a video at different bitrates and divides it into multiple segments. To achieve a high quality of experience (QoE), the clients using the HAS determine the video bitrate of the upcoming segments according to the network

conditions [3]–[5]. Commercialized HAS solutions include Microsoft's smooth streaming, Apple's HTTP live streaming, Adobe's HTTP dynamic streaming, and international standard named dynamic adaptive streaming over HTTP (DASH) [6]–[9].

Several types of researches have been carried out for adaptive streaming schemes in recent years [10]–[13]. The majority of the proposed adaptive streaming schemes determine the video bitrate of the upcoming segments on the client-side. Each client determines the video bitrate of upcoming segments based on the estimated bandwidth and requests the segments periodically. When the multiple clients stream a video over the same access network, the network bandwidth can be estimated inaccurately due to the segment request pattern of clients. The inaccurately estimated bandwidth causes unnecessary bitrate switching and video stalling events, leading to the degradation in the QoE of clients [14], [15]. These problems lead to inefficient resource utilization and unfair bitrate allocation in view of multiple clients. In mobile

networks where the channel conditions change dynamically, the degradation in the QoE of clients, the inefficient resource utilization, and the unfair bitrate allocation become even worse. The fundamental reason behind the aforementioned problems is that the clients are oblivious to the radio channel conditions and cannot coordinate with each other in adaptive streaming.

Multi-access edge computing (MEC) can be adopted to address the problems caused when the multiple clients stream a video over the same access network [16]. By deploying servers in the radio access network (RAN), the MEC brings the capabilities for computation and storage to the edge of mobile networks. In adaptive streaming, the MEC server can play a key role through the central management for network resources and connected clients by utilizing the information of RAN and application in real-time. The clients can coordinate with each other and determine the video bitrate appropriate to the channel conditions with the assistance of the MEC server. Therefore, edge computing-assisted adaptive streaming is able to jointly optimize the QoE of clients, network resource utilization, and fairness.

There are key challenges for edge computing-assisted adaptive streaming. First, the adaptive streaming framework should be designed to be compatible with the DASH standard. The framework should support the resource management of networks and coordination among clients without any modification in both the servers and clients. Next, the bitrate adaptation should jointly consider QoE, resource utilization, and fairness. These factors correlate with each other. Providing high video bitrate to the clients can improve the QoE and resource utilization but lead to unfairness among clients. Furthermore, the video bitrate of the specific client may be sacrificed to improve fairness among clients in some cases. The clients need to switch the video bitrate for fair bitrate allocation but the frequent and abrupt bitrate switching degrades the QoE of clients.

In this paper, we propose an edge computing-assisted adaptive streaming scheme for mobile networks. We first present an edge computing-assisted adaptive streaming framework compatible with the DASH standard. Then, we design an optimization model that jointly considers the QoE and fairness. The joint optimization model can balance the QoE and fairness while satisfying the constraints for network resources and clients. We present a greedy-based bitrate allocation algorithm to efficiently solve the joint optimization problem formulated on the basis of the designed model. Finally, we evaluate the performance of the proposed scheme through extensive simulation based on mobile network configurations. The experimental results show that the proposed scheme outperforms the existing schemes in terms of the QoE, resource utilization, and fairness.

The rest of the paper is organized as follows. We discuss related work in Section II and illustrate the proposed scheme in Section III. The simulation-based evaluation for the proposed scheme is presented in Section IV, and we conclude the paper in Section V.

## II. RELATED WORK

Several adaptive streaming schemes to improve the QoE of mobile clients have been proposed in recent years. Mekki *et al.* have proposed an adaptive streaming scheme using indoors-outdoors detection [17]. In the HAS, the clients use the segment throughput to determine the video bitrate of the upcoming segments. However, due to user mobility and channel fading, the bitrate adaptation only using segment throughput is not effective in mobile video streaming. The proposed scheme detects the indoors-outdoors states based on the received signal strength of clients. The video bitrate of the upcoming segments is determined by considering the indoors-outdoors states and the buffer status of clients. Bethanabhotla *et al.* have proposed an adaptive streaming scheme to improve the QoE of clients in wireless networks where multiple clients compete for limited bandwidth [18]. The proposed scheme targets on designing an optimal scheduling policy for video-on-demand (VoD) services over wireless networks. To determine the optimal video bitrate for multiple clients, the proposed scheme formulates a network utility maximization (NUM) problem. The network utility includes the requested video bitrate, stability in client buffer, and fairness among clients. The Lyapunov optimization method is adopted for the proposed scheme to efficiently solve the network utility maximization problem. Based on the optimization results, the bitrate adaptation helpers determine the video bitrate for multiple clients. Xie *et al.* have proposed the energy-efficient cache resource allocation and QoE optimization for the HAS over cellular networks [19]. The main objective of the proposed scheme is to maximize the QoE of clients as well as the energy cost saving in the base stations. The high video bitrate can enhance the QoE of clients but the energy cost also can increase due to the high transmission delay. To jointly optimize the QoE and energy consumption, the proposed scheme formulates a content cache management problem for adaptive streaming. The proposed scheme solves the content cache management problem by dividing it into two sub-problems. The content to cache on the base stations is determined according to the joint optimization results.

The existing adaptive streaming schemes for multi-client mobile networks largely rely on the client-driven approach. These schemes utilize cross-layer information and various optimization strategies to determine the video bitrate for multiple clients. However, the existing schemes are still lack of coordination among clients and achieve suboptimal performance in highly fluctuating environment such as mobile networks. To remedy this, the server and network-assisted DASH (SAND-DASH) can be adopted [20]. The key idea of the SAND-DASH is that the servers and in-network elements help the bitrate adaptation through the collaboration with the clients. Li *et al.* have proposed a wireless bottleneck coordination scheme using the SAND-DASH framework [21]. The proposed scheme performs the optimization for the resource-constrained utility to improve the QoE of clients in a shared wireless link. The optimal coordination

points are determined by using a resource pricing solution. Cofano *et al.* have performed the design and experimental evaluation for the network-assisted HAS strategies [22]. The performance of the two approaches is evaluated to find the optimal solution for adaptive streaming based on software-defined networking (SDN). The first approach is to allocate the network bandwidth slices to the video flows, and the second approach is to perform the bitrate guidance for the clients. The evaluation results show that the network-assisted strategies can improve the QoE of clients and allocate the video bitrate among the clients fairly while utilizing the available network resources efficiently. Bentaleb *et al.* have proposed an SDN-based resource management architecture for the HAS [23]. The main objective of the proposed architecture is to achieve the QoE maximization and fair bitrate allocation. The network resources can be managed in programmatic and flexible manners by using the capabilities of SDN. The proposed architecture determines the network resources to be allocated to each client based on the expected QoE of clients.

However, the additional modifications in the servers, clients, and in-network elements are required to support the SAND-DASH in adaptive streaming scheme, and it is hard to access mobile network information in real-time with SAND-DASH assistance. Besides the approaches based on the SAND-DASH, adaptive streaming schemes with edge computing have been studied in recent years. These schemes do not demand the additional modifications in the servers and clients, and the network resources and clients can be managed by accessing mobile network information in real-time. Wang *et al.* have proposed an adaptive wireless video streaming scheme based on edge computing [24]. The proposed scheme presents an adaptive transcoding framework with the capabilities of edge computing. The transcoding servers are deployed close to the base stations hence can utilize the information of channel and clients in the adaptive transcoding. The proposed scheme formulates the adaptive transcoding strategy as the NUM problem and divides it into two sub-problems. Based on the optimization results, the proposed scheme allocates the coding resources to the clients adaptively and manages the available network resources. Mehrabi *et al.* have proposed an edge computing-assisted adaptive mobile video streaming scheme [25]. The proposed scheme designs the optimization model jointly considering the QoE of clients, network resource utilization, and fairness among clients. Moreover, the proposed scheme presents the method of distributing the connected clients among multiple edge servers. To efficiently solve the clients to edge servers mapping and the bitrate selection problem, the proposed scheme presents a near-optimal greedy-based scheduling algorithm.

Yan *et al.* have proposed a hybrid edge cloud and client adaptation scheme for the HAS in cellular networks [26]. The edge cloud performs the joint optimization of QoE and fairness by using the information of RAN and application. The QoE model is designed in the form of a continuum to estimate the cumulative viewing experience of clients during the video streaming sessions. The QoE continuum model jointly considers the requested video bitrate, video stalling events, and bitrate switching. The proposed scheme presents a heuristic algorithm to efficiently solve the joint optimization problem derived from the QoE continuum model. Rahman *et al.* have proposed an edge computing-assisted joint quality adaptation scheme for mobile video streaming [27]. The proposed scheme presents a joint throughput estimation method based on the capabilities of edge computing. Using the results of the joint throughput estimation, the edge cloud can assign the video bitrate among clients fairly, and the unnecessary bitrate switching can be reduced. The proposed scheme designs an integer non-linear programming (INLP) model to jointly optimize the QoE of competing clients in mobile networks. Moreover, a heuristic bitrate selection algorithm is presented to efficiently solve the joint optimization problem.

The edge computing-assisted adaptive streaming scheme should satisfy the following requirements to achieve the optimal performance. First, the adaptive streaming framework should be compatible with the DASH standard. The proposed scheme in this paper shifts the adaptation intelligence from the clients to the edge cloud which includes the base stations of mobile networks and the MEC servers. The edge cloud considers the constraints for network resources and clients to determine the video bitrate for multiple clients. Moreover, any modification in both the HAS server and clients is not needed since the edge cloud can capture the segment requests of clients. The next requirement is that the joint optimization problem should be designed by considering the QoE of clients, network resource utilization, and fairness among clients. However, it is hard to achieve the optimal QoE, resource utilization, and fairness simultaneously since these factors affect each other in the joint optimization. From this motivation, the proposed scheme designs the joint optimization model that can balance the QoE and fairness. The proposed scheme then presents the greedy-based bitrate allocation algorithm to efficiently solve the joint optimization problem. The proposed scheme can determine the video bitrate maximizing the QoE and resource utilization and minimizing the loss in fairness through the bitrate allocation algorithm.

## III. PROPOSED SCHEME
### A. EDGE COMPUTING ASSISTED ADAPTIVE STREAMING FRAMEWORK

Fig. 1 shows the edge computing-assisted adaptive streaming framework for mobile networks. The HAS server stores a video by encoding into the different bitrates and dividing it into the multiple segments of equal playback length. The HAS server is connected to the edge cloud through the content delivery network (CDN). The edge cloud consists of the base station and the MEC server. Thanks to the proximity to the base station and the management capability of edge computing, the MEC server can access the information of RAN and application in real-time. The clients connected to
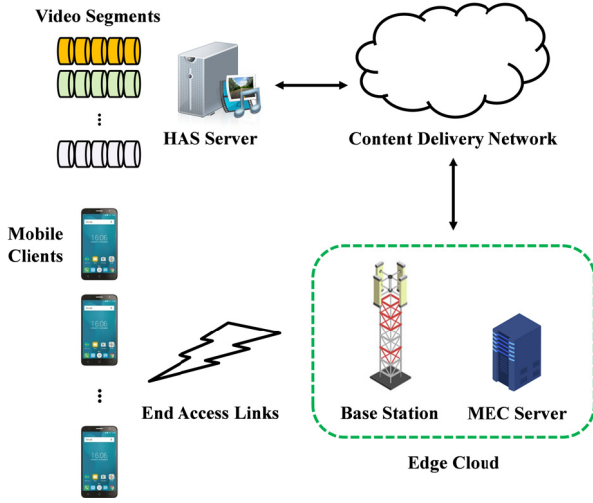
**FIGURE 1.** Edge computing-assisted adaptive streaming framework for mobile networks.

the edge cloud compete for the available bandwidth of the end access links and use the conventional HAS to determine the video bitrate of the upcoming segments. During the streaming sessions, the clients report the HAS information such as the requested video bitrate and the buffer status periodically to the edge cloud. This process is feasible because the 3GPP has standardized the reporting process for the QoE metrics based on the HTTP POST protocol [28]. To determine the video bitrate for multiple clients, the MEC server performs the joint optimization by considering the information of RAN and application. The clients request the segments periodically based on the HAS mechanism. The MEC server captures the segment requests of clients and modifies those requests by newly determining the video bitrates based on the joint optimization.

We assume that all the clients stream the same video and the number of clients does not change during the streaming sessions. The latency due to the capturing of the segment requests is ignored in the joint optimization. Moreover, the device characteristics of clients such as resolution, battery, and computation capability are homogeneous. The proposed scheme does not impose additional modifications to the radio resource scheduling process of the base stations. The base stations allocate the available resource blocks (RBs) to the clients in a proportionally fair (PF) manner to support the clients efficiently [29]. In other words, the available RBs are allocated to the clients by considering the channel quality between the clients and the base stations and the amount of data communicated in the clients.

From the next sub-section, we illustrate the key components of the proposed scheme. The list of the parameters involved in the proposed scheme and their descriptions are summarized in Table 1.

## B. CHANNEL ESTIMATION
It is important to estimate the transmission rate of channel accurately for the channel-aware bitrate adaptation. Both the

**TABLE 1.** Description of parameters involved in the proposed scheme.

| Notation | Description |
|---|---|
| $M$ | Number of clients connected to the edge cloud |
| $r_{max}, r_{min} \in R$ | Maximum and minimum video bitrate in bitrate set $R$ |
| $Thr_{i,n}, Thr_{i,n}^s$ | Instant and smoothed channel throughput for the client $i$ at the $n$th reporting period |
| $S_{i,n}$ | Amount of received data reported from the client $i$ at the $n$th reporting period |
| $T_{period}$ | Time cycle of the reporting process |
| $\alpha$ | Weighting parameter for channel estimation |
| $C_i$ | Number of requested segments for the client $i$ up to the current |
| $r_{i,k}$ | Video bitrate of the $k$th segment requested by the client $i$ |
| $AB_i, AS_i$ | Accumulated video bitrate and bitrate switching for the client $i$ |
| $UF_i$ | Accumulated unfairness level for the client $i$ |
| $\beta$ | Weighting parameter to give the penalty to the bitrate switching in the joint optimization |
| $W_k$ | Number of available resource blocks at the base stations when requesting the $k$th segment |
| $B_{i,k}$ | Buffer level for the client $i$ after receiving the $k$th segment |
| $B_{max}$ | Maximum buffer level of clients |
| $r_{max}^{thr}, r_{max}^{buf}$ | Sustainable bitrates relying on segment throughput and client buffer level |
| $\delta_S, \delta_{UF}$ | Switching and unfairness thresholds to control the bitrate switching and fairness in the greedy-based bitrate allocation algorithm |

base stations and the clients perform various strategies to manage the available network resources efficiently in mobile networks. The base stations change the modulation coding scheme (MCS) according to the channel status to minimize the waste of available network resources. The clients periodically report the channel-related metrics such as channel quality indicator (CQI) and signal-to-interference-noise ratio (SINR) to the base stations. Moreover, the clients report the amount of actually received data to the base stations. The base stations allocate the available RBs to the clients by considering the feedback information of clients. To estimate the effective bandwidth of the transmission channel, the instant channel throughput is calculated by using the amount of received data and the reporting period.

$$Thr_{i,n} = \frac{S_{i,n}}{T_{period}} \qquad (1)$$

$Thr_{i,n}$ and $S_{i,n}$ are the instant channel throughput and the amount of received data for the client $i$ measured at the $n$th reporting period, respectively. $T_{period}$ is the time cycle that the clients report the channel-related metrics to the base stations. Typically, the time cycle of the reporting process ranges from tens to hundreds of milliseconds. We set the value of the time cycle as a fixed value in the proposed scheme. The transmission rate of channel fluctuates severely even during a short time due to user mobility and channel fading. Therefore, the transmission rate of the channel can be estimated

inaccurately when only using the instant channel throughput. We apply the exponentially weighted moving average (EWMA) to reduce the impact of the channel fluctuations on the channel estimation. The smoothed channel throughput is calculated by using the weighting parameter and the instant channel throughput.

$$Thr_{i,n}^s = \alpha \cdot Thr_{i,n-1}^s + (1 - \alpha) \cdot Thr_{i,n} \tag{2}$$

$Thr_{i,n}^s$ is the smoothed channel throughput for the client $i$ measured at the $n$ th reporting period. $\alpha$ is the weighting parameter to smooth the instant channel throughput, which can improve the robustness of the channel estimation for the channel fluctuations. We set the value of the weighting parameter as a fixed value in the proposed scheme.

### C. QUALITY OF EXPERIENCE AND FAIRNESS

According to several studies on the QoE of the HAS, the video bitrate, bitrate switching, video stalling, and startup delay jointly affect the QoE of clients. However, compared with the other influencing factors, the impact of startup delay on the QoE of clients is not significant. The main influencing factors should be considered for efficient QoE optimization. To this end, the proposed scheme formulates the joint optimization problem by considering the video bitrate, bitrate switching, and video stalling.

Delivering high video bitrate to the clients is important to improve the QoE of clients. If the playback bitrate is high, the clients can view a video clearly. However, providing high video bitrate to the clients does not always improve the QoE of clients since there is a tradeoff between the video bitrate and the video stalling. The probability of video stalling increases when the video bitrate is high because the available network bandwidth is not always sufficient to support the high bitrate. To consider the impact of video bitrate on the QoE of clients during the streaming sessions, the accumulated video bitrate is calculated by aggregating the requested bitrate for the clients.

$$AB_i = \sum_{k=1}^{C_i} r_{i,k} \tag{3}$$

$AB_i$ is the video bitrate accumulated for the client $i$ up to the current. $C_i$ is the index of segment currently requested at the client $i$, and $r_{i,k}$ is the video bitrate of the $k$ th segment requested at the client $i$. The joint optimization in the proposed scheme targets achieving the high value of the accumulated video bitrate without video stalling events. Therefore, the clients can utilize the available network resources efficiently during the streaming sessions while ensuring the smooth video playback.

Switching the video bitrate of clients is a key function of the HAS. The clients can experience the smooth video playback by matching the video bitrate to the network conditions. However, the frequent and abrupt bitrate switching degrades the QoE of clients. The proposed scheme uses the difference between the video bitrates of the consecutive segments requested by the clients as the QoE metric for the bitrate

switching. To consider the impact of bitrate switching on the QoE of clients during the streaming sessions, the accumulated bitrate switching is calculated by aggregating the bitrate difference for the clients.

$$AS_i = \sum_{k=2}^{C_i} |r_{i,k} - r_{i,k-1}| \tag{4}$$

$AS_i$ is the bitrate difference accumulated for the client $i$ up to the current. The proposed scheme regards the accumulated bitrate switching as a penalty in the joint optimization to minimize the unnecessary bitrate switching.

Due to the prominent role of video stalling in determining the QoE of clients, avoiding stalling events is important during the streaming sessions. The stalling events occur when the client buffer is entirely consumed. After receiving the video segments, the proposed scheme checks whether the stalling events occur by using the client buffer level calculated on the basis of the smoothed channel throughput and the playback length of segments.

$$B_{i,k} = B_{i,k-1} - \left( \frac{r_{i,k} \cdot \tau}{Thr_{i,n}^s} \right) + \tau \tag{5}$$

$B_{i,k}$ is the buffer level for the client $i$ after receiving the $k$ th segment of the video bitrate $r_{i,k}$, and $\tau$ is the playback length of segments. The proposed scheme regards the video stalling as a constraint in the joint optimization to avoid stalling events during the streaming sessions.

The base stations of mobile networks usually schedule the available resources to the connected clients according to the PF policy. The number of RBs to be allocated to the clients is determined by accounting for the channel quality and the amount of data communicated. The clients select the best sustainable bitrate for the upcoming segments according to the resource share determined by the PF scheduler at the base stations. The proposed scheme considers the average of bitrates allocated to the other clients for the fair bitrate allocation among clients. The average of bitrates allocated to the other clients is calculated to measure the fairness level of clients.

$$\bar{r}_k = \frac{1}{M-1} \sum_{j=1, j \neq i}^{M} r_{j,k} \tag{6}$$

$\bar{r}_k$ is the average of bitrates allocated to the other clients for the $k$ th segment. $M$ is the number of clients connected to the edge cloud. The objective of the fair bitrate allocation is to minimize the bitrate deviations among clients during the streaming sessions. The accumulated unfairness level of clients is calculated as the sum of the differences between the average of bitrates allocated to the other clients and the bitrate selected by the clients.

$$UF_i = \sum_{k=1}^{C_i} |\bar{r}_k - r_{i,k}| \tag{7}$$

$UF_i$ is the unfairness level accumulated for the client $i$ up to the current. The proposed scheme regards the accumulated unfairness level as a penalty in the joint optimization to minimize the bitrate deviations among clients. However, due to

the correlation between QoE and fairness, some clients may sacrifice the video bitrate to achieve the fair bitrate allocation. Switching the video bitrate of clients is required to improve the fairness among clients but the QoE of clients may be damaged by the frequent and abrupt bitrate switching. In the next sub-section, we illustrate how the joint optimization problem is formulated by considering the correlation between QoE and fairness.

### D. JOINT OPTIMIZATION PROBLEM

Note that the influencing factors to the QoE of clients have different characteristics with each other. The proposed scheme combines the accumulated video bitrate, bitrate switching, and unfairness level by using the weighting parameters to represent these factors as the one utility value. In particular, the weighting parameter for the unfairness level of clients is used to balance the QoE and fairness. The weighting parameter to balance the QoE and fairness is calculated on the basis of the previous unfairness level of clients.

$$\rho_i = 1 - \frac{\left| r_{i,k-1} - \bar{r}_{k-1} \right|}{r_{max} - r_{min}} \quad (8)$$

$\rho_i$ is the weighting parameter to balance the QoE and fairness. $r_{max}$ and $r_{min}$ are the maximum and minimum video bitrate, respectively. When the previous unfairness level of clients is high, we need to focus on improving the fairness among clients. In the opposite case, we need to focus on improving the QoE of clients instead of allocating the video bitrate among clients fairly.

The joint optimization for multiple clients is formulated as a utility maximization problem that considers the correlation between QoE and fairness, and the joint optimization problem includes the following integer non-linear constraints.

$$\text{Maximize } U_i = \rho_i \cdot (AB_i - \beta \cdot AS_i) - (1 - \rho_i) \cdot UF_i \quad (9)$$

$$\text{Subject to } \sum_{i=1}^{M} \left\lceil \frac{r_{i,k}}{Thr_{i,n}^s} \right\rceil \leq W_k \quad (10)$$

$$0 < B_{i,k} \leq B_{max} \quad (11)$$

$$r_{i,k} \in R \quad (12)$$

The objective function (9) aims to maximize the utility value of clients while balancing the QoE and fairness. $\beta$ is the weighting parameter to give the penalty to the bitrate switching. The constraint (10) means that the total resources distributed for the clients should not exceed the available resources of the base stations. The joint optimization problem specifies the constraint (11) to ensure the video playback without interruptions during the streaming sessions. The client buffer level should be larger than zero and lower than the maximum buffer level to avoid the buffer underflow and overflow, respectively. The constraint (12) means that the video bitrate allocated to the clients should belong to the set of video bitrates available at the HAS server.

The joint optimization problem in (9)-(12) is hard to solve within the polynomial-time due to the existence of the integer non-linear constraints. Brute-force search can be used to investigate all the possibilities of determining the video bitrate with the maximum achievable utility value to solve the joint optimization problem. However, the complexity of the brute-force search grows exponentially with the increase in the number of clients and constraints. In the next sub-section, we illustrate the bitrate allocation algorithm that determines the video bitrate for multiple clients in a greedy manner and satisfies the objectives of the joint optimization accordingly.

### E. ONLINE OPTIMIZATION ALGORITHM

The details of the greedy-based bitrate allocation algorithm are described in Algorithm 1. For the first video segment, all the clients determine the video bitrate based on the conventional HAS. The bitrate allocation algorithm is activated after all the clients receive the first segment. When the clients request the video segments to the HAS server, the video bitrate that jointly optimizes the QoE of clients, resource utilization, fairness among clients is determined through the bitrate allocation algorithm. The segment request information of the clients is modified according to the results of the bitrate allocation.

The basic idea of the greedy-based bitrate allocation algorithm is to determine the video bitrate jointly maximizing the resource utilization and the utility value of clients. The bitrate allocation algorithm first seeks the maximum available bitrate lower than the sustainable bitrates while satisfying the resource constraint. The sustainable bitrate relying on the segment throughput, $r_{max}^{thr}$, is the maximum bitrate which is lower than the segment throughput. The sustainable bitrate relying on the client buffer level, $r_{max}^{buf}$, is the maximum bitrate which does not cause the stalling events. The bitrate allocation algorithm compares the difference between the bitrate selected in the aforementioned step and the previously requested bitrate with the switching threshold. The switching threshold is calculated by using the sustainable bitrate relying on the client buffer level.

$$\delta_S = \left| r_{i,k-1} - r_{max}^{buf} \right| \quad (13)$$

$\delta_S$ is the threshold to control the bitrate switching of clients in the bitrate allocation algorithm. The clients can utilize the available resources efficiently without experiencing the unnecessary bitrate switching when the bitrate difference is lower than the switching threshold. At the next step, the bitrate allocation algorithm considers three cases associated with fairness among clients.

If the bitrate satisfying the switching constraint is higher than the previously requested bitrate, the bitrate allocation algorithm sets the unfairness threshold, $\delta_{UF}$, to the fixed value of 0.5. Then, the bitrate allocation algorithm calculates the unfairness level of clients, $|r - \bar{r}_{k-1}| / (r_{max} - r_{min})$, which takes a value between 0 and 1. The video bitrate of the upcoming segments is determined by comparing the unfairness threshold with the unfairness level of clients. As a result, the QoE of clients can be improved with high video bitrate and less bitrate switching while maintaining a certain level of fairness.

---

**Algorithm 1** Greedy-Based Bitrate Allocation

$r_{i,k}^{avg}$: average bitrate for the client $i$ up to the $k$th segment

$Thr_{i,k}^{seg}$: segment throughput for the client $i$ after receiving the $k$th segment

1:     **Compute** the sustainable bitrates $r_{max}^{thr}$ and $r_{max}^{buf}$

2:     **Compute** the switching and unfairness thresholds $\delta_S$ and $\delta_{UF}$

3:     **For** each bitrate $r \in R$ in decreasing order

4:       **If** allocation of $r$ satisfies the resource constraint

          **AND** $r \le max(r_{max}^{thr}, r_{max}^{buf})$

5:         **If** $|r - r_{i,k-1}| \le \delta_S$

6:           **If** $r > r_{i,k-1}$

7:             **Determine** $\delta_{UF}$ to the fixed value

8:             **If** $|r - \bar{r}_{k-1}| / (r_{max} - r_{min}) \le \delta_{UF}$

9:               $r_{i,k} = r$; **Break;**

10:           **If** $r == r_{i,k-1}$ **AND** $|r - \bar{r}_{k-1}| / (r_{max} - r_{min}) \le \delta_{UF}$

11:             $r_{i,k} = r$; **Break;**

12:           **If** $r < r_{i,k-1}$

13:             **If** $r > r_{i,k-1}^{avg}$

14:               $r_{i,k} = r$; **Break;**

15:             **Else**

16:               $r_{i,k} = r_{i,k-1}$; **Break;**

17:     **If** $r_{i,k} == 0$

18:       **For** each bitrate $r \in R$ in decreasing order

19:         **If** allocation of $r$ satisfies the resource constraint **AND** $r \le max(r_{max}^{thr}, r_{max}^{buf})$

20:           $r_{i,k} = r$; **Break;**

21:     **If** $r_{i,k} == 0$

22:       $r_{i,k} = max\{r | r < Thr_{i,k-1}^{seg}\}$; **Break;**

23:     **Update** auto-tuning parameter $\rho_i$

24:     **Compute** $AB_i$, $AS_i$, $UF_i$, and $U_i$

25:     **Update** $B_{i,k}$

26:     **Return** $U_i$

---

If the bitrate satisfying the switching constraint is equal to the previously requested bitrate, the bitrate allocation algorithm focuses on balancing the QoE and fairness. Minimizing bitrate switching improves the QoE of clients but the bitrate deviations among clients may increase. The bitrate allocation algorithm in this case determines the unfairness threshold by using the sustainable bitrates relying on segment throughput and client buffer level. The average sustainable bitrate is first calculated to determine the unfairness threshold.

$$r_{sus}^{avg} = \frac{\left(r_{max}^{thr} + r_{max}^{buf}\right)}{2} \qquad (14)$$

$r_{sus}^{avg}$ is the video bitrate that is able to utilize the available network resources efficiently and ensure the smooth video playback of clients. The bitrate allocation algorithm calculates the unfairness threshold by using the average sustainable bitrate.

$$\delta_{UF} = \frac{\left|r_{sus}^{avg} - \bar{r}_{k-1}\right|}{r_{max} - r_{min}} \qquad (15)$$

When the average sustainable bitrate is high, the bitrate allocation algorithm provides the maximum available bitrate to the clients depending on the average of bitrates allocated to the other clients. The bitrate allocation algorithm in the opposite case strives to improve fairness among clients while minimizing the degradation in QoE.

If the bitrate satisfying the switching constraint is lower than the previously requested bitrate, the bitrate allocation algorithm uses the average bitrate of clients to determine the video bitrate of the upcoming segments. When the bitrate satisfying the switching constraint is higher than the average bitrate of clients, the bitrate allocation algorithm determines the bitrate satisfying the switching constraint as the video bitrate of the upcoming segments. In the opposite case, the bitrate allocation algorithm does not change the video bitrate to minimize the unnecessary bitrate switching and the degradation in the average bitrate.

If there is no such bitrate available, the bitrate allocation algorithm looks for the maximum bitrate which satisfies the resource constraint and does not experience the stalling events. If still there is no such bitrate available, the bitrate allocation algorithm then selects the maximum bitrate lower than the segment throughput. After the video bitrate of the upcoming segments is determined, the bitrate allocation algorithm updates the weighting parameter to balance the QoE and fairness. Furthermore, the client buffer level is updated according to the determined bitrate in the joint optimization. The bitrate allocation algorithm finally returns the utility value of clients which is obtained by using the objective function (9).

## IV. PERFORMANCE EVALUATION

### A. SIMULATION SETUP

To verify the performance of the proposed scheme, we have built an ns-3 based simulation environment that includes the HAS server, the clients, and the edge cloud [30]. The HAS server stores a video by encoding into 10 bitrate levels [184, 414, 644, 874, 1104, 1334, 1564, 2024, 2484, and 2944 kbps]. The number of clients connected to the edge cloud is set to 20. The playback length of segments and the maximum buffer size of clients are set to 2 s and 30 s, respectively. We set the weighting parameter $\alpha$ and $\beta$ used to the channel estimation and the joint optimization model to 0.8 and 1.0, respectively. The clients request the video segments based on the HAS. The simulation is conducted for 3600 s in the performance evaluation.

For the network setup, we configure the LTE network environment and consider a single mobile cell in which the base station and the clients are uniformly distributed. We set the time cycle of the reporting process used in the channel estimation to 250 ms. The clients start the streaming session sequentially and move within a cell at a fixed speed. In other

**TABLE 2.** Mobile network configuration.

| Parameter | Corresponding Value |
|---|---|
| Number of UEs | 20 |
| UE Distribution | Uniform |
| Path Loss Model | Hata Model PCS Extension |
| BS Transmission Power | 38 dBm |
| UE Distance | 250 ~ 500 m |
| Channel Bandwidth | 5 MHz |
| Number of Downlink RBs | 25 |
| Scheduler | Proportional Fairness |
| UE Speed | 3 km/h (Urban) and 60 km/h (Vehicular) |

words, the clients experience different channel conditions. During the streaming session, the channel throughput of clients varies according to the client's moving speed. The list of the mobile network configuration parameters and their corresponding values are summarized in Table 2.

We compare the performance of the proposed scheme with the existing edge computing-assisted adaptive streaming schemes. In the performance evaluation, we refer to these schemes proposed in [25] and [26] as ECAA and Prius, respectively. The ECAA is an edge computing-assisted adaptive mobile video streaming scheme that includes the clients to the edge servers mapping and the joint optimization of QoE and fairness. We adopt the ECAA for a single-cell scenario with one edge server. The Prius is a hybrid edge cloud and client adaptation scheme for cellular networks. Using the QoE continuum model, the Prius jointly optimizes the QoE and fairness. We use the parameter settings of the existing schemes described in their papers. We take the average value of the five experiments in the performance evaluation. For the performance evaluation of the proposed scheme over time-varying mobile networks, we define the fading scenarios as Urban and Vehicular according to the client's moving speed. In the Urban scenario where the clients are moving slowly, the fluctuations of channel throughput are low and the average channel throughput is high. In the Vehicular scenario where the clients are moving fastly, due to the channel fading, the fluctuations of channel throughput are high and the average channel throughput is low.

### B. EVALUATION METRICS

We utilize the QoE metrics described in [31] for the performance comparison between the proposed scheme and the existing schemes. The average bitrate is calculated by using the per-segment bitrate requested at the clients during the streaming sessions.

$$r_i^{avg} = \frac{1}{N_i} \sum_{k=1}^{N_i} r_{i,k} \qquad (16)$$

$r_i^{avg}$ is the average bitrate for the client $i$ during the streaming sessions, and $N_i$ is the number of requested video segments for the client $i$ during the streaming sessions.

The switching frequency of clients denotes how frequent the video bitrate is changed during the streaming sessions. We calculate the switching frequency of clients by accumulating the number of changes in the bitrate level of clients. Moreover, we calculate the switching magnitude of clients to measure how abrupt the video bitrate is changed during the streaming sessions.

$$r_i^{fre} = \sum_{k=2}^{N_i} f\left(r_{i,k-1}, r_{i,k}\right) \qquad (17)$$

$$r_i^{mag} = \frac{1}{N_i - 1} \sum_{k=1}^{N_i-1} \left| r_{i,k+1} - r_{i,k} \right| \qquad (18)$$

$r_i^{fre}$ is the switching frequency for the client $i$ during the streaming sessions. $f(r_{i,k-1}, r_{i,k})$ is the indicator function to check the changes in the bitrate level of clients. If the values of $r_{i,k-1}$ and $r_{i,k}$ are different, the value of the function is set to 1, and it is set to zero when the values of $r_{i,k-1}$ and $r_{i,k}$ are the same. $r_i^{mag}$ is the switching magnitude for the client $i$ during the streaming sessions.

The resource utilization is also a key metric to evaluate the performance of the proposed scheme. To efficiently utilize the available resources of networks, the clients should request the video segments with the bitrate fit to the channel throughput. The resource utilization of clients, $\varphi_{res}^{avg}$, is calculated by using the ratio between the requested video bitrate and the channel throughput.

$$\varphi_{res}^{avg} = \frac{1}{N_i} \sum_{k=1}^{N_i} \left( \frac{r_{i,k}}{Thr_{i,n}^S} \right) \qquad (19)$$

We use the Jain's fairness index to measure the fairness level of the clients [32]. The Jain's fairness index, $J_F$, is calculated by using the average bitrate of competing clients.

$$J_F = \frac{\left( \sum_i r_i^{avg} \right)^2}{M \cdot \sum_i \left( r_i^{avg} \right)^2} \qquad (20)$$

Before taking the average value of the five experiments in the performance evaluation, the average bitrate, switching frequency, switching magnitude, and resource utilization are averaged by the number of clients.

### C. RESULTS

Fig. 2 shows the average bitrate for the compared schemes according to the fading scenarios. As shown in the results, the proposed scheme achieves the highest average bitrate in all the fading scenarios. The reason is that the proposed scheme minimizes the degradation in the average bitrate of clients through the greedy-based bitrate allocation algorithm. The Prius jointly optimizes the QoE and fairness by adjusting the QoE continuum of clients. However, the video bitrate of clients may be sacrificed to achieve the fair QoE continuum among clients. The ECAA determines the video bitrate of the upcoming segments through the greedy-based bitrate selection algorithm. Using the switching and fairness thresholds, the greedy-based bitrate selection algorithm jointly optimizes the QoE and fairness. In the ECAA, the switching threshold
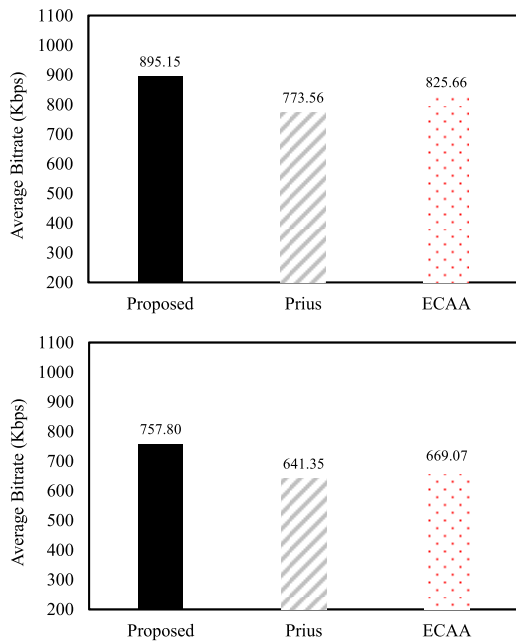
**FIGURE 2.** Average bitrate for the compared schemes according to the fading scenarios (top: Urban, bottom: Vehicular).
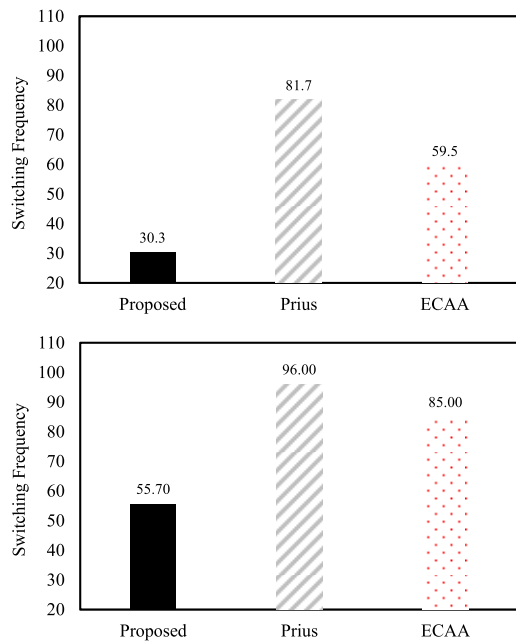


**FIGURE 3.** Switching frequency for the compared schemes according to the fading scenarios (top: Urban, bottom: Vehicular).



**FIGURE 4.** Switching magnitude for the compared schemes according to the fading scenarios (top: Urban, bottom: Vehicular).

**TABLE 3.** Resource utilization and Jain's fairness index according to the fading scenarios (top: Urban, bottom: Vehicular).

| Metrics | Proposed | Prius | ECAA |
|---|---|---|---|
| Resource Utilization | 0.93 | 0.85 | 0.84 |
| Jain's Fairness Index | 0.87 | 0.83 | 0.84 |

| Metrics | Proposed | Prius | ECAA |
|---|---|---|---|
| Resource Utilization | 0.90 | 0.83 | 0.80 |
| Jain's Fairness Index | 0.85 | 0.83 | 0.81 |

can be wrongly determined according to the channel fluctuations, and the fairness threshold is fixed for all the clients. The average bitrate of clients degrades since the video bitrate is changed unnecessarily and sacrificed to achieve the fair bitrate allocation.

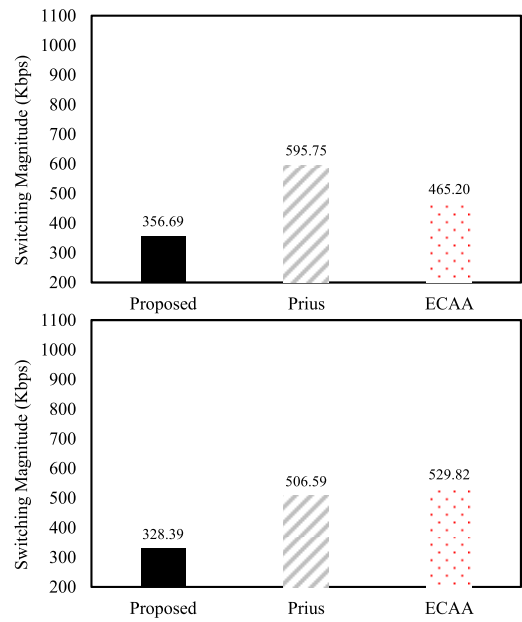Fig. 3 and 4 show the switching frequency and magnitude for the compared schemes according to the fading scenarios, respectively. The proposed scheme achieves the lowest switching frequency and magnitude in all the fading scenarios. The reason is that the greedy-based bitrate allocation algorithm in the proposed scheme balances the QoE and fairne ss, minimizing the unnecessary bitrate switching. On the other hand, the Prius and ECAA suffer from the frequent and abrupt bitrate switching during the streaming sessions. To achieve the fair QoE continuum among clients, the video bitrate of clients may be changed unnecessarily in the Prius. Due to the switching threshold wrongly determined when the channel highly fluctuates, the ECAA may change the video bitrate of clients unnecessarily.

We confirmed that all the compared schemes did not experience the video stalling events during the streaming sessions. The proposed scheme and ECAA check whether the clients experience the video stalling events after receiving the video segments. Moreover, if there is no such bitrate available, the proposed scheme and ECAA determine the maximum sustainable bitrate for the clients to ensure the smooth video playback. The Prius regards the video stalling events as the penalty in the joint optimization. To improve the QoE continuum among clients, the Prius determines the maximum

sustainable bitrate for the clients if there is no such bitrate available.

The resource utilization and Jain's fairness index for the compared schemes according to the fading scenarios are summarized in Table 3. As shown in the results, the proposed scheme achieves the highest resource utilization in all the fading scenarios and minimizes the loss in fairness. The reason is that the greedy-based bitrate allocation algorithm in the proposed scheme maximizes the resource utilization of clients while balancing the QoE and fairness.

## V. CONCLUSION

In this paper, we have proposed the edge computing-assisted adaptive streaming scheme for mobile networks. Due to the lack of channel awareness and coordination among clients, the existing client-driven approaches suffer from the degradation in QoE, inefficient resource utilization, and unfair bitrate allocation. With the assistance of edge computing, the proposed scheme shifts the adaptation intelligence of clients to the edge cloud. To optimize the QoE of clients, resource utilization, and fairness among clients, the proposed scheme performs the joint optimization by utilizing the information of RAN and application. The joint optimization problem is formulated by considering the correlation between QoE and fairness. The proposed scheme presents the greedy-based bitrate allocation algorithm to efficiently solve the joint optimization problem. Through the performance evaluation, we have proved that the proposed scheme can improve the QoE of clients and resource utilization even in the fluctuating channel conditions and minimize the loss in fairness.

## REFERENCES

[1] Cisco. *Cisco Visual Networking Index: Forecast and Trends 2018–2023 White Paper*. Accessed: Mar. 2020. [Online]. Available: https://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/white-paper-c11-741490.html.

[2] A. Bentaleb, B. Taani, A. C. Begen, C. Timmerer, and R. Zimmermann, "A survey on bitrate adaptation schemes for streaming media over HTTP," *IEEE Commun. Surveys Tuts.*, vol. 21, no. 1, pp. 562–585, 1st Quart., 2019.

[3] A. A. Barakabitze, N. Barman, A. Ahmad, S. Zadtootaghaj, L. Sun, M. G. Martini, and L. Atzori, "QoE management of multimedia streaming services in future networks: A tutorial and survey," *IEEE Commun. Surveys Tuts.*, vol. 22, no. 1, pp. 526–565, 1st Quart., 2020.

[4] M. T. Vega, C. Perra, F. De Turck, and A. Liotta, "A review of predictive quality of experience management in video streaming services," *IEEE Trans. Broadcast.*, vol. 64, no. 2, pp. 432–445, Jun. 2018.

[5] S. Petrangeli, J. V. D. Hooft, T. Wauters, and F. D. Turck, "Quality of experience-centric management of adaptive video streaming services: Status and challenges," *ACM Trans. Multimedia Comput., Commun., Appl.*, vol. 14, no. 2s, pp. 1–29, May 2018.

[6] Microsoft. *Microsoft Smooth Streaming Specifications*. Accessed: Jan. 2018. [Online]. Available: https://azure.microsoft.com/en-us/services/media-services/

[7] Apple. *HTTP Live Streaming (HLS) Draft Specification*. Accessed: Dec. 2017. [Online]. Available: http://tools.ietf.org/html/draft-pantos-http-live-streaming-12

[8] Adobe. *Adobe HTTP Dynamic Streaming (HDS) and Media Manifest (F4M) Specifications*. Accessed: Dec. 2017. [Online]. Available: https://www.adobe.com/devnet/hds.html

[9] MPEG. *Dynamic Adaptive Streaming Over HTTP (DASH) Specification*. Accessed: Jan. 2018. [Online]. Available: https://mpeg.chiariglione.org/standards/mpeg-dash

[10] C. Liu, I. Bouazizi, and M. Gabbouj, "Rate adaptation for adaptive HTTP streaming," in *Proc. ACM Conf. Multimedia Syst. (MMSys)*, Feb. 2011, pp. 169–174.

[11] T.-Y. Huang, R. Johari, N. McKeown, M. Trunnell, and M. Watson, "A buffer-based approach to rate adaptation: Evidence from a large video streaming service," in *Proc. ACM Conf. Special Interest Group Data Commun. (SIGCOMM)*, Aug. 2014, pp. 187–198.

[12] J. Jiang, V. Sekar, and H. Zhang, "Improving fairness, efficiency, and stability in HTTP-based adaptive video streaming with festive," in *Proc. Int. Conf. Emerg. Netw. Exp. Technol. (CoNEXT)*, Dec. 2012, pp. 97–108.

[13] Z. Li, X. Zhu, J. Gahm, R. Pan, H. Hu, A. C. Begen, and D. Oran, "Probe and adapt: Rate adaptation for HTTP video streaming at scale," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 4, pp. 719–733, Apr. 2014.

[14] J. Chen, R. Mahindra, M. A. Khojastepour, S. Rangarajan, and M. Chiang, "A scheduling framework for adaptive video delivery over cellular networks," in *Proc. 19th Annu. Int. Conf. Mobile Comput. Netw. (MobiCom)*, 2013, pp. 389–400.

[15] D. D. Vleeschauwer, H. Viswanathan, A. Beek, S. Benno, G. Li, and R. Miller, "Optimization of HTTP adaptive streaming over mobile cellular networks," in *Proc. IEEE Conf. Comput. Commun. (INFOCOM)*, Apr. 2013, pp. 989–997.

[16] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 4, pp. 2322–2358, Aug. 2017.

[17] S. Mekki, T. Karagkioules, and S. Valentin, "HTTP adaptive streaming with indoors-outdoors detection in mobile networks," in *Proc. IEEE Conf. Comput. Commun. Workshops (INFOCOM WKSHPS)*, May 2017, pp. 671–676.

[18] D. Bethanabhotla, G. Caire, and M. J. Neely, "Adaptive video streaming for wireless networks with multiple users and helpers," *IEEE Trans. Commun.*, vol. 63, no. 1, pp. 268–285, Jan. 2015.

[19] J. Xie, R. Xie, T. Huang, J. Liu, and Y. Liu, "Energy-efficient cache resource allocation and QoE optimization for HTTP adaptive bit rate streaming over cellular networks," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2017, pp. 1–6.

[20] *Dynamic Adaptive Streaming Over HTTP (DASH) Part 5: Server and Network Assisted DASH (SAND)*, document ISO/IEC 23009-5. Accessed: Feb. 2017. [Online]. Available: https://www.iso.org/standard/69079.html

[21] Z. Li, S. Zhao, D. Medhi, and I. Bouazizi, "Wireless video traffic bottleneck coordination with a DASH SAND framework," in *Proc. Vis. Commun. Image Process. (VCIP)*, Nov. 2016, pp. 1–4.

[22] G. Cofano, L. De Cicco, T. Zinner, A. Nguyen-Ngoc, P. Tran-Gia, and S. Mascolo, "Design and experimental evaluation of network-assisted strategies for HTTP adaptive streaming," in *Proc. 7th Int. Conf. Multimedia Syst.*, May 2016, pp. 1–12.

[23] A. Bentaleb, A. C. Begen, and R. Zimmermann, "SDNDASH: Improving QoE of HTTP adaptive streaming using software defined networking," in *Proc. 24th ACM Int. Conf. Multimedia*, Oct. 2016, pp. 1296–1305.

[24] D. Wang, Y. Peng, X. Ma, W. Ding, H. Jiang, F. Chen, and J. Liu, "Adaptive wireless video streaming based on edge computing: Opportunities and approaches," *IEEE Trans. Services Comput.*, vol. 12, no. 5, pp. 685–697, Sep. 2019.

[25] A. Mehrabi, M. Siekkinen, and A. Y-Jääski, "Edge computing assisted adaptive mobile video streaming," *IEEE Trans. Mobile Comput.*, vol. 18, no. 4, pp. 787–800, Apr. 2019.

[26] Z. Yan, J. Xue, and C. W. Chen, "Prius: Hybrid edge cloud and client adaptation for HTTP adaptive streaming in cellular networks," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 27, no. 1, pp. 209–222, Jan. 2017.

[27] W. U. Rahman, C. S. Hong, and E.-N. Huh, "Edge computing assisted joint quality adaptation for mobile video streaming," *IEEE Access*, vol. 7, pp. 129082–129094, 2019.

[28] *Progressive Download and Dynamic Adaptive Streaming Over HTTP*, document TS 26.247 V12.1.0, 3GPP. Accessed: Dec. 2013. [Online]. Available: http://goo.gl/4EJbvd

[29] L. Li, M. Pal, and Y. R. Yang, "Proportional fairness in multi-rate wireless LANs," in *Proc. IEEE 27th Conf. Comput. Commun. (INFOCOM)*, Apr. 2008, pp. 1678–1686.

[30] NSNAM. *The Network Simulator NS-3*. Accessed: Aug. 2019. [Online]. Available: https://nsnam.org

[31] X. Yin, A. Jindal, V. Sekar, and B. Sinopoli, "A control-theoretic approach for dynamic adaptive video streaming over HTTP," in *Proc. ACM Conf. Special Interest Group Data Commun. (SIGCOMM)*, Aug. 2015, pp. 325–338.

[32] A. B. Sediq, R. H. Gohary, R. Schoenen, and H. Yanikomeroglu, "Optimal tradeoff between sum-rate efficiency and Jain's fairness index in resource allocation," *IEEE Trans. Wireless Commun.*, vol. 12, no. 7, pp. 3496–3509, Jul. 2013.

**MINSU KIM** received the B.S. degree from the Electronics and Communications Engineering Department, Kwangwoon University, Seoul, South Korea, in 2017, where he is currently pursuing the Ph.D. degree. His research interests include QoS, QoE support, multimedia systems, and streaming protocols.

**KWANGSUE CHUNG** (Senior Member, IEEE) received the B.S. degree from the Electrical Engineering Department, Hanyang University, Seoul, South Korea, the M.S. degree from the Electrical Engineering Department, Korea Advanced Institute of Science and Technology (KAIST), Seoul, and the Ph.D. degree from the Electrical Engineering Department, University of Florida, Gainesville, FL, USA. Before joining the Kwangwoon University, in 1993, he spent ten years with the Electronics and Telecommunications Research Institute (ETRI), as a Member of the Research Staff. From 1991 to 1992, he was an Adjunct Professor with KAIST. From 2003 to 2004, he was a Visiting Scholar with the University of California at Irvine, Irvine, CA, USA. His research interests include communication protocols and networks, QoS mechanisms, and video streaming.

• • •