

# **Resource-Efficient QoS-Aware Video Streaming Using UAV-Assisted Networks**

Downloaded from: https://research.chalmers.se, 2024-05-09 08:19 UTC

Citation for the original published paper (version of record):

Bhar, C., Ghosh, D., Agrell, E. (2024). Resource-Efficient QoS-Aware Video Streaming Using UAV-Assisted Networks. IEEE Transactions on Cognitive Communications and Networking, 10(2): 649-659. http://dx.doi.org/10.1109/TCCN.2023.3336908

N.B. When citing this work, cite the original published paper.

© 2024 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, or reuse of any copyrighted component of this work in other works.

This document was downloaded from http://research.chalmers.se, where it is available in accordance with the IEEE PSPB Operations Manual, amended 19 Nov. 2010, Sec, 8.1.9. (http://www.ieee.org/documents/opsmanual.pdf).

# Resource-efficient QoS-Aware Video Streaming using UAV-Assisted Networks

Chayan Bhar, Member, IEEE, Debayani Ghosh, Erik Agrell, Fellow, IEEE

Abstract-Emerging video services are associated with stringent quality-of-service (QoS) and high data-rate requirements. Moreover, the presence of data-rate-hungry mobile users in future networks necessitate sophisticated design strategies. The deployment of unmanned aerial vehicle access point (UAP)assisted networks (UANs) has been proposed to ensure high data-rates to mobile users. Moreover, UAPs can be equipped with energy-efficient caches to facilitate video delivery with stringent QoS. However, the mobility of users and UAPs may cause temporal variations in the QoS experienced by users. This paper conducts an extensive performance evaluation of a UAN, by studying the effect of user behavior, mobility of users and UAPs, and a temporal variation of video popularity on the QoS. The QoS is measured in terms of the delay experienced by the users. To that end, a time-dependent queueing model and its associated fluid approximation models are derived, which are illustrated to be reasonably accurate in an appropriate asymptotic regime. A detailed analysis of these models reveals that low delay, i.e., high QoS, can be ensured in UANs. Finally, a reinforcement-learning (RL) approach based on these models is utilized to minimize the number of deployed UAPs and the playout buffer size while guaranteeing a certain QoS.

*Index Terms*—UAV-assisted networks, QoS-aware video delivery, time-dependent queueing model.

#### I. INTRODUCTION

There has been a rapid increase in data traffic generated by mobile video streaming services in recent years. Internet video is expected to contribute up to 82% of all Internet traffic in 2022 [1]. Video traffic is also expected to increase with the increasing use of high resolution video formats. However, such traffic can congest existing networks in the near future due to their bandwidth-intensive nature. Network congestion will in turn affect the quality-of-service (QoS) [2]. However, overprovisioning the network infrastructure to accommodate the growing video traffic can result in high energy consumption and an expensive transport network. Furthermore, existing and future networks comprise highly mobile users that can cause a spatio-temporal variation in the user density and video request patterns. This necessitates migration to new networks, like Unmanned aerial vehicle access point-assisted networks (UANs) that can offer stringent QoS [1] with efficient usage of

E. Agrell is with the Department of Electrical Engineering, Chalmers University of Technology, SE-41296 Gothenburg, Sweden.

the access points and the playout buffers, which in this paper are called "resource efficiency" for short.

UANs employing mobile unmanned aerial vehicle access points (UAPs) have been proposed to provide consistent communication links to mobile end users in next generation networks [3]. The UAPs can be co-deployed with the existing base stations, whenever there is a sudden increase in the number of mobile users. Moreover, UAPs can be strategically deployed to follow mobile users on a need basis in geographic regions with high user density or data requirements [4]–[6]. This makes UAPs more effective and cost-efficient compared to deploying static ultra-dense small cells [7]. Therefore, UAPs provide the flexibility of managing a time-varying density of users [3], [8].

On the other hand, edge caching is proposed as a potential solution to the bandwidth congestion problem in core and metro networks [9]. Edge caching can be implemented in UANs by equipping UAPs with light and energy-efficient small fog caches that can cache popular bandwidth-intensive on-demand videos [10]. Such UAPs can stream popular videos to mobile users with extremely high data rates [11], [12]. Therefore, such a scheme can ensure stringent QoS and network bandwidth efficiency, while simultaneously decreasing the bandwidth consumption of the transport network, as videos are streamed from network locations near the end users.

#### A. Literature Survey

We divide the related literature in three categories: (i) UAP placement and trajectory optimization for its power efficiency, (ii) user-dependent UAP deployment strategies, and (iii) cache placement and video caching strategies in UANs. Trajectory optimization of a unmanned aerial vehicle (UAV) using various algorithms is illustrated to maximize the throughput and power efficiency of an individual UAP in [5], [6], [13], [14]. Power-efficient allocation of a UAP from a set of deployed UAPs that have fixed locations for a particular coverage area using optimal transport theory and user behavior prediction is discussed in [8] and [12], respectively. Semi-definite relaxation and coordinate descent methods are used to design a bandwidth- and power-efficient scheme in which UAPs equipped with caches serve users having heterogeneous datarate requirements [10]. Placement of a UAP to maximize its coverage in a scenario with QoS bounds using the exhaustive search method is discussed in [15].

The deployment and placement of a UAP as a function of the users' parameters like (i) user density and mobility [3]–[5], [11], [16], (ii) end user throughput [6], (iii) video

This work was funded in parts by SERB under grant SRG/2020/001364 and by the Swedish Research Council (VR) under grant 2021-03709.

C. Bhar is with the Department of Electronics and Communication Engineering, National Institute of Technology Warangal, Telangana, India. (e-mail: cbhar@nitw.ac.in).

D. Ghosh is with the Department of Electronics and Communication Engineering, Thapar Institute of Engineering and Technology, Patiala, Punjab, India.

requests [12], and (iv) QoS [15], [17]-[19] requirements have also been studied in the literature. It is illustrated that mobile users experience poor QoS, i.e., high delay if the time to deploy a UAP is too high [5]. High data-rates resulting in good QoS can be achieved by increasing the number of deployed caches in small-cell stationary access points using a greedy method [20]. Optimal cache placement, power allocation, and trajectory optimization subject to the mobility constraints of a UAP is illustrated to maximize user throughput using the block-coordinate descent and successive convex approximation methods [21]. UAP deployment as a function of the time-averaged request arrival and download rates is illustrated to minimize the video streaming delay in [18], [22]. User mobility and energy-efficiency-aware UAP movement are illustrated to limit delay [17] and the average transmit power [19]. Finally, caching of contents in a high-altitude platform using federated learning is implemented for estimating content popularity in [23]. Optimal cache placement using the entropy weight method is illustrated to minimize the content access delay in a scenario with multiple access points (APs) and cache units [24]. The delay experienced by users at concurrent time intervals is used for designing the transmission power of UAVs in [12]. Delay can be experienced by a user due to (i) content retrieval from the caches [21]-[24], (ii) request queueing at UAPs, and (iii) video streaming from a UAP [5], [12], [20], [22].

Reinforcement-learning (RL)-based strategies for energy minimization in UAV-assisted wireless powered sensor networks and for learning the optimal control policy to minimize delay in queueing networks have been proposed in [25], [26], respectively. However, such methods have not been implemented for designing UAN deployment as a function of QoS and user mobility parameters.

#### B. Motivation

The deployment of UANs has been proposed for edge caching and last mile content delivery using UAVs thereby reducing backhaul traffic and improving QoS [12], [27], [28]. Even a modest cache size of 100 Mbits can reduce the backhaul traffic by more than 50% in favorable cases [29, Fig. 4]. Mobile UAPs can be deployed in or traverse through different geographical locations in which they can experience a spatio-temporal variation in the video popularity [12], [30]. Similarly, the simultaneous mobility of UAPs and users can cause loss of connectivity and a variation in the download times observed by users [16]. Inconsistent download rates while streaming videos from the UAPs can result in poor QoS. Moreover, UAPs suffer from energy, space, and size limitations, which limit the hardware capability and equipment complexity of UAPs [3], [5], [8], [9], [13]. Therefore, the maximum user density and frequency of video requests that can be served by a UAP is limited and affected by its energy [15]. For example, if a UAP is low on charge, it can stream to a small number of users [8] or stream videos with low popularity [31], [7]. The number of users that can be served by UAPs in a UAN depends on the deployment of UAPs and is timedependent. Thus, the deployment of UAPs must be managed



Fig. 1: Architecture for the proposed UAN.

in real time as a function of the videos' popularity and size, user mobility pattern [7], and the constraints related to UAP deployment [12]. Otherwise, the efficiency and effectiveness of UANs and the real-time QoS experienced by users can be affected. Thus, QoS is an important parameter to decide the optimal number of UAPs required.

UAPs can be deployed to serve users assuming that the QoS bounds allow some waiting time to the network provider before starting to stream a video. However, the user parameters required for UAP deployment in a geographical region are not available to individual UAPs. Furthermore, the energy limitation in UAPs discourages on-board implementation of computationally complex deployment strategies. However, existing literature does not discuss QoS-aware deployment of UAPs as functions of the time varying request pattern, user arrival, and channel conditions. Furthermore, minimizing the resources, i.e., number of UAPs deployed and the playout buffer size is also absent in existing literature.

#### C. Our contributions

Whereas previous literature on cache-enabled UANs has considered system parameters such as video popularity, user mobility, and channel conditions to be constant, this paper takes one step further and investigates how the experienced QoS (queueing delay) is affected under the more realistic scenario of time-varying system parameters. Using an analytic model and extensive simulations, it is shown that the temporal variations of these parameters, not only their averages, significantly affect the experienced QoS. This allows real-time UAP deployment for a resource-efficient UAN design.

The rest of the paper is arranged as follows. In the next section we describe the considered UAN and provide a detailed description of the system model. This is followed by an analysis of the derived model, results on QoS performance of the proposed UAN, and design of UAP deployment strategies using RL in Section III. Finally the paper concludes in Section IV with a discussion and conclusion.

# II. SYSTEM MODEL

The considered UAN architecture is illustrated in Fig. 1. It is similar to the architecture in [9], [24] but with coordinated



Fig. 2: Network of queues for the UAN.

UAPs. It is also similar to the architecture in [3], [31], which we extend by considering a time-varying deployment of UAPs. In this section, we formulate an analytical model for the system of Fig. 1. The symbols used for developing the queueing model are described in Table I while those used for developing the user mobility, channel, and video popularity models are described in Table II.

## A. Functions of the central processing unit (CPU) and UAPs

The proposed UAN simultaneously employs traditional base stations and UAPs to provide consistent network throughput to all users. Each UAP connects to the CPU C via base stations using fronthaul links. The CPU C is centrally located and performs processing of the payload data from UAPs and base stations, and video popularity estimation. This allows UAPs and base stations to carry minimal processing elements for energy efficiency. Thus, in our proposed scheme, the UAPs are used only for content delivery. All baseband functionalities are executed in C, which also connects UAPs and base stations to the core network. For this purpose, at any given time, each UAP is connected to C through a base station using wireless fronthaul links. UAPs forward the channel state information and video requests from users to the CPU over fronthaul links using the scheme discussed in [32]. Thus, at any given point of time, the CPU has knowledge of the channel state information and video request arrival from all UAPs. The CPU informs UAPs about video streaming requirements. UAPs position themselves in the two-dimensional space with a fixed height to serve the users mandated by C while conforming to their energy limitations. The UAPs prevent high interference to users in their respective coverage area by using the reinforcement learning-based mobility strategy of [33]. UAPs are shown to cause minimum interference to users when they hover below 10 m [34]. Thus, in most of our simulated scenarios we fix the UAP height to 10 m. Moreover, the UAPs utilize the interference-minimization techniques described in [35] to maximize end-user data rates. The CPU C utilizes an estimated streaming rate for all users throughout the UAN. It also utilizes UAPs to facilitate fast streaming of bandwidth-intensive videos, thereby providing high QoS at spatio-temporal locations that have a high intensity of video

TABLE I: List of symbols used for the queueing model. All time-varying functions are defined for time t > -T, unless otherwise stated.

Symbol	Description
Ι	Number of UAPs deployed
$D_i$	$i^{th}$ UAP consisting of FC $F_i$
C	CPU
v	Considered video
$\kappa$	Size of v
$\zeta(t)$	The maximum number of simultaneous $v$ streams supported by $C$ , i.e., UAPs deployed at $t$ , upper and lower limits $\zeta_{\rm h}$ , $\zeta_{\rm l}$ , period of deployment $T_{\zeta}$
$\lambda(t)$	Ensemble mean of $v$ request inter-arrival time to $C$ at $t$
r(t)	Ensemble mean of the rate of download at $t$
$\phi$	Probability of users to disconnect from $C$ without downloading $v$
$\gamma$	Mean time after which requests from $Q(t)$ retry $v$ download, desired $\gamma_{\rm M}$
$\alpha$	Mean time after which requests are abandoned while waiting for $v$ in $W(t)$
au, s	Running variables of t. The "test" request arrives at $t = \tau$
W(t)	Number of requests enqueued in the service queue at C for requests of v, scaled queue $W^{\eta}(t)$ , fluid limit $W^{f}(t)$ , diffusion limit $W^{d}(t)$
$\hat{W}^{\mathrm{f}}(t),$ $\hat{W}^{\mathrm{d}}(t)$	Fluid and diffusion limits of $W(t)$ for $t > \tau$
O(t)	Number of requests enqueued in the retrial queue
	for users that left $W(t)$ and can request $v$ from $C$ , scaled queue $Q^{\eta}(t)$ , fluid limit $Q^{f}(t)$ , diffusion
$\mathbf{\Omega}(t)$	The system of queues consisting of $W(t)$ and $Q(t)$ considered for analysis, scaled system $\Omega^{\eta}(t)$ , fluid
A(t)	imit $\mathbf{M}^{r}(0)$ , diffusion limit $\mathbf{M}^{r}(t)$ Request arrival process for $v$ at $C$ , scaled process
$\Delta(t)$	Departure process due to completion of v stream- ing from C, scaled process $\Delta^{\eta}(t)$ , fluid limit $\Delta^{f}(t)$ , diffusion limit $\Delta^{d}(t)$
$E^{\mathrm{f}}(t)$	Fluid limit of service initiation process at UAPs with derivative $e^{f}(t)$ , diffusion limit $E^{d}(t)$ , and scaled process $E^{\eta}(t)$
$P^{\mathrm{f}}(t)$	Fluid limit of attainment process, scaled process $P^{\eta}(t)$
-T	Operator for processes observed at $t \ge \tau$ Lower limit of time for observing all processes at $C$ to calculate delay
$\delta^{\mathrm{f}}(t)$	Fluid limit of delay in downloading v, scaled delay $\delta^{\eta}(t)$ diffusion limit $\delta^{d}(t)$ desired $\delta_{V}$
$X_i(t)$	Intensity for a time-inhomogeneous Poisson process $\Pi_i()$

requests [31]. The UAPs operate as cache helpers in a femtocaching environment [20]. For this purpose, each UAP  $D_i$ ,  $i \in \{1, \ldots, I\}$  is equipped with a small, energy-, weight-, and space-efficient femto cache (FC)  $F_i$  [5], [31], while the CPU Cis equipped with a cloud unit. The FCs are significantly lightweight compared to the cloud unit. The decision to stream a video v from the FC of a UAP to a requesting user is taken by C on a per-request basis.

TABLE II: List of symbols used for the user mobility, channel, and video popularity models and the resource optimization problem.

Symbol	Description
$\langle \Delta r \rangle$	User displacement from initial position
$\alpha_1, \alpha_2$	Parameters of distribution for $\Delta t$ , $\Delta r$
S(t)	number of distinct locations visited by a user till $t$
$\rho, \nu$	parameter of user mobility, predicting a user's
	tendency to explore a new location
$\alpha_3$	Path loss exponent
$f_{\rm c}, c, B$	Carrier frequency, speed of light, and bandwidth
	available to a UAP-user link
h, d(t)	height of UAP, distance of user from UAP on
	ground plane
$\delta_1,  \delta_2$	Excessive path loss exponent for line-of-sight and
	non-line-of-sight links between UAP and user
$p_{\text{LoS}}$ ,	probability of line-of-sight and non-line-of-sight
$p_{ m NLoS}$	links
$\sigma^2$	AWGN noise power
$\psi, \psi_1$	Environment-dependent constant for modeling the
	wireless channel
p(t),	UAP-to-user transmit power, path loss between
P(t)	UAP and user
$P_i(t)$	Path loss experienced by interefering signals from
	UAP $i$ to a particular user
$\alpha_4, \alpha_5$	Parameters for direct and word-of-mouth recom-
	mendation of $v$
q	Intrinsic popularity of the video
N(t)	Number of users in a geographical region, lower
	limit $N_{\rm l}$ , upper limit $N_{\rm h}$ , period of variation $T_N$
$\beta_1$ ,	Constants of the optimization problem in Section
$a_1, \cdots, a_5$	III-B
Ξ	Total interference experienced by a user from all
	UAPs

### B. Network model

The model formulated in this paper is with respect to a particular video v, CPU C, and UAPs  $D_i$ ,  $i \in \{1, \ldots, I\}$ . It is assumed that physical resources of UAPs like storage capacity, battery, energy efficiency, cache placement, etc., are managed according to [5]. The interplay between different videos due to the size limitation of FCs is not considered for simplicity of analysis. Since video popularity varies spatiotemporally, a CPU can experience a temporal variation in its observed video popularity [30]. The popularity is modeled by the inter-arrival time between requests for a video v. Empirical evidence suggests that video popularity is modeled by a Zipf distribution. However, for a single video on a short timescale, the request arrival process can be modelled by a pseudostationary Poisson point process [36]. Thus, the request interarrival time, a measure of video popularity, is exponentially distributed with ensemble mean  $\lambda(t)$  at time t [30], [36], [37]. When a user requests v, the request is transferred to C, which employs UAPs to deliver v to requesting users. A network operator owns the terrestrial network and C and pays the mobile network operators owning  $D_i$  for streaming videos. The business model for the UAN can be similar to [38], [39]. For deciding the payment for streaming videos, the network operator considers the ensemble mean of inter-arrival time  $\lambda(t)$  and the QoS experienced by end users. Poor channel conditions, low battery charge, or low  $\lambda(t)$  can require a high

payment.

A UAP  $D_i$  caches v depending on its popularity [30]. [36] and the payment available from C [40]. Therefore, a measure of video popularity is required both at the CPU and by mobile network operators owning  $D_i$ . The CPU shares the estimated  $\lambda(t)$  with the mobile network operators to decide video caching in  $D_i$ . When a user requests a video v, C coordinates its delivery from the FC of a UAP connected to Cand having v [40]. The v streaming rate r(t) depends on the propagation channel between the UAP and user and the battery state of the concerned UAP [16], [41]. The download time for v with size  $\kappa$  is assumed to be exponentially distributed with mean  $\kappa/r(t)$  for all users. As UAPs can be mobile, the users downloading v can experience a temporal variation in the streaming rates and hence in the download data-rates. The ensemble mean of the download time is assumed to be equal for all users for simplicity of analysis. The total number of users that can simultaneously stream v from UAPs associated with C is limited to  $\zeta(t)$ . The CPU controls  $\zeta(t)$  and I by varying the payments provided to UAPs.

At any point in time, if  $\zeta(t)$  UAPs simultaneously stream v to an equal number of users, then further requests for vwait in a buffer located at C for completion of the ongoing streams. A user with a queued request may lose interest in v, resulting in deletion of the queued request. Moreover, the link quality between a user with a queued request and the CPU may momentarily deteriorate due to network outage resulting from high interference, airflow disturbances, beam steering errors, etc. [42]. Such events result in deletion of the queued requests after an exponentially distributed interval with mean  $\alpha$  [43]. Therefore, high  $\alpha$  can arise from frequent channel disruptions or high user mobility. If the link with a user deteriorates, then the request is dropped with probability  $\phi$  and reenqueued otherwise. If it is reenqueued, it happens after an exponentially distributed interval with mean  $\gamma$  [16], [44]–[46]. If v corresponds to a video segment in a fractional caching scheme [47], [48],  $\gamma$  can be proportional to the time for emptying the playout buffer [11]. Therefore, if users are equipped with small playout buffers, then  $\gamma$  must be small.

#### C. Queueing model

Below we formulate a queueing model for the CPU C using the assumptions of Section II-A and the queue illustration of Fig. 2. The model consists of a service queue and a retrial queue. The service queue W(t) comprises the users downloading v from UAPs and the users for which the corresponding requests are queued at C. The retrial queue Q(t) consists of users that requested v but their requests were removed from W(t) due to channel deterioration. Thus, W(t) and Q(t) are positive integers that represent lengths of the service and retrial queues, respectively. Requests from Q(t) are again placed in W(t) after  $\gamma$  s. We assume that W(t) is a first-in-first-out type of queue. If there is an exogenous request for v at time t and the number of users streaming v at t is less than  $\zeta(t)$ , then  $W(t) < \zeta(t)$ . In this case, v is streamed immediately by a UAP to the requesting user. In contrast, if  $\zeta(t)$  users are simultaneously streaming v at time t then  $W(t) = \zeta(t)$  and the exogenous request is placed in W(t).  $\gamma$  and  $\zeta(t)$  are upper limited to  $\gamma_{\rm M}$  and  $\zeta_{\rm h}$ , respectively. Deletion of queued requests from W(t) is with rate  $\alpha$  due to the reasons mentioned in Section II-A. The network resource optimization performed in this paper is with respect to  $\gamma$  and  $\zeta(t)$  as described in Section III-B using (34) and (35).

In the context of the proposed UAN, the arrival and service processes in the queueing system of Fig. 2 are time-dependent Markovian processes. Moreover, the number of simultaneous v streams from C is also time-dependent. Therefore, the queueing model for C is an  $M(t)/M(t)/\zeta(t)$  queue, in which W(t) allows abandonment, while Q(t) facilitates retrial [49]. In general, such a queueing model is analytically intractable [49]. However, to make this model amenable to analysis, a fluid approximation is often adopted in the queuing literature. In the asymptotic regime wherein the arrival rate and the number of simultaneous streams from the service queue are large, the fluid approximation has been shown to be accurate. Motivated by this, the proposed system in this work is modeled along similar lines as in [49].

We denote with  $\Pi_i(X_i(t))$  independent, inhomogeneous Poisson point processes with time-varying intensity functions  $X_i(t)$  [50, Sec. 5.1]. The length of the service queue at time t increases with the number of v requests (i) present initially, W(0), (ii) arriving from the retrial queue,  $(1/\gamma)\Pi_1(Q(t))$ , and (iii) arriving exogenously,  $\Pi_2(1/\lambda(t))$ . On the other hand, the service queue decreases with the number of v requests (i) leaving the service queue,  $(1/\alpha)\Pi_3((W(t) - \zeta(t))^+)$ where  $x^+ = \max(x, 0)$  and (ii) completing v streaming  $(1/\kappa)\Pi_4(\min(W(t), \zeta(t)))r(t)$ .

The length of the retrial queue at time t increases with the number of v requests present initially at Q(0) and the requests routed to Q(t) from W(t). It decreases with the enqueueing of  $(1/\gamma) \int_0^t \Pi_1(Q(s)) ds$  requests, from Q(t) to W(t). Therefore, the sample paths [51, Sec. 5.2] for W(t) and Q(t) are the unique solutions of the equations

$$W(t) = W(0) + \frac{1}{\gamma} \int_0^t \Pi_1 (Q(s)) \, \mathrm{d}s + \int_0^t \Pi_2 \left(\frac{1}{\lambda(s)}\right) \, \mathrm{d}s \\ - \frac{1}{\alpha} \int_0^t \Pi_3 \left( (W(s) - \zeta(s))^+ \right) \mathrm{d}s \\ - \frac{1}{\kappa} \int_0^t \Pi_4 \left( \min \left( W(s), \zeta(s) \right) \right) r(s) \, \mathrm{d}s, \tag{1}$$

$$Q(t) = Q(0) + \frac{(1-\phi)}{\alpha} \int_0^t \Pi_3 \left( (W(s) - \zeta(s))^+ \right) ds - \frac{1}{\gamma} \int_0^t \Pi_1 \left( Q(s) \right) ds$$
(2)

where  $1/\lambda(t)$ , r(t), and  $\zeta(t)$  are locally integrable. We define  $\Omega(t) = (W(t), Q(t))$ . Using the theory of strong approximations for Poisson processes [49], a random sample path construction of C is performed to do an asymptotic sample path analysis and obtain the fluid limit theorems. In the asymptotic regime,  $\zeta(t)$  is scaled up in response to a similar scaling up of the arrival rate by customers. The asymptotic regime

in *C* denoted by  $\eta$  corresponds to  $\Omega^{\eta}(t) = (W^{\eta}(t), Q^{\eta}(t))$ , where  $\Omega^{\eta}(t)$  is the fluid approximation of  $\Omega(t)$ . The scaled parameters are the initial conditions  $\Omega^{\eta}(0) = \lceil \eta \Omega^{f}(0) + \sqrt{\eta} \Omega^{d}(0) \rceil + o(\sqrt{\eta})$  for constants  $\Omega^{f}(0) = (W^{f}(0), Q^{f}(0))$ and  $\Omega^{d}(0) = (W^{d}(0), Q^{d}(0)), \eta/\lambda(t)$ , and  $\eta\zeta(t)$  [49]. The scaled processes are

$$W^{\eta}(t)$$

$$= W^{\eta}(0) + \frac{1}{\gamma} \int_{0}^{t} \Pi_{1} \left( Q^{\eta}(s) \right) \, \mathrm{d}s + \int_{0}^{t} \eta \Pi_{2} \left( \frac{1}{\lambda(s)} \right) \, \mathrm{d}s \\ - \frac{1}{\alpha} \int_{0}^{t} \Pi_{3} \left( \left( W^{\eta}(s) - \eta \zeta(s) \right)^{+} \right) \mathrm{d}s \\ - \frac{1}{\kappa} \int_{0}^{t} \Pi_{4} \left( \min \left( W^{\eta}(s), \eta \zeta(s) \right) \right) r(s) \, \mathrm{d}s,$$
(3)

$$Q^{\eta}(t)$$

$$= Q^{\eta}(0) + \frac{(1-\phi)}{\alpha} \int_{0}^{t} \Pi_{3} \left( (W^{\eta}(s) - \eta\zeta(s))^{+} \right) \mathrm{d}s$$
$$- \frac{1}{\gamma} \int_{0}^{t} \Pi_{1} \left( Q^{\eta}(s) \right) \, \mathrm{d}s.$$
(4)

To obtain the fluid approximation, an asymptotic regime is considered wherein the arrival rate  $1/\lambda(t)$  and  $\zeta(t)$  are scaled up by a factor  $\eta > 0$ . Hence, in this new system  $(W^{\eta}(t), Q^{\eta}(t))$ , the arrival and service rates are  $\eta/\lambda(t)$  and  $\eta r(t)/$ , respectively [52, Chapter 6]. Using the strong law of large numbers [51, App. A.2], the fluid limits are obtained for all t > 0 as

$$\lim_{\eta \to \infty} \frac{1}{\eta} W^{\eta}(t) = W^{\mathrm{f}}(t), \quad \lim_{\eta \to \infty} \frac{1}{\eta} Q^{\eta}(t) = Q^{\mathrm{f}}(t) \tag{5}$$

with the initial conditions

$$\lim_{\eta \to \infty} \frac{1}{\eta} W^{\eta}(0) = W^{f}(0), \quad \lim_{\eta \to \infty} \frac{1}{\eta} Q^{\eta}(0) = Q^{f}(0).$$
(6)

The convergences of (5) are uniform on any compact subset of  $t \ge 0$  [49], while the convergences at time t = 0 follow from the strong law of large numbers [51, App. A.2] [53]. The validity of this assumption is illustrated in Section III.

The asymptotic assumptions made in (5) and (6) allow us to derive  $W^{f}(t)$  and  $Q^{f}(t)$  by dividing (3) and (4) by  $\eta$  on both sides and taking the limits as  $\eta \to \infty$ 

$$W^{f}(t) = W^{f}(0) + \frac{1}{\gamma} \int_{0}^{t} \Pi_{1} \left( Q^{f}(s) \right) ds + \int_{0}^{t} \Pi_{2} \left( \frac{1}{\lambda(s)} \right) ds - \frac{1}{\alpha} \int_{0}^{t} \Pi_{3} \left( \left( W^{f}(s) - \zeta(s) \right)^{+} \right) ds - \frac{1}{\kappa} \int_{0}^{t} \Pi_{4} \left( \int_{0}^{t} \min \left( W^{f}(s), \zeta(s) \right) \right) r(s) ds, \quad (7)$$

$$Q^{f}(t) = Q^{f}(0) + \frac{(1-\phi)}{\alpha} \int_{0}^{t} \Pi_{3} \left( \left( W^{f}(s) - \zeta(s) \right)^{+} \right) ds - \frac{1}{\gamma} \int_{0}^{t} \Pi_{1} \left( Q^{f}(s) \right) ds.$$

$$(8)$$

Finally, differentiating (7) and (8) with respect to t, we obtain the rate of change in the two fluid limits  $W^{f}(t)$  and  $Q^{f}(t)$ , which are given by the solutions of

$$\frac{\mathrm{d}}{\mathrm{d}t}W^{\mathrm{f}}(t) = \frac{1}{\lambda(t)} + \frac{Q^{\mathrm{f}}(t)}{\gamma} - \frac{r(t)\min\left(W^{\mathrm{f}}(t),\zeta(t)\right)}{\kappa} - \frac{\left(W^{\mathrm{f}}(t) - \zeta(t)\right)^{+}}{\alpha}, \qquad (9)$$

$$\frac{\mathrm{d}}{\mathrm{d}t}Q^{\mathrm{f}}(t) = \frac{(1-\phi)\left(W^{\mathrm{f}}(t)-\zeta(t)\right)^{+}}{\alpha} - \frac{Q^{\mathrm{f}}(t)}{\gamma}.$$
 (10)

The queue length  $W^{\rm f}(t)$  increases due to the arrival of exogeneous and retrial requests for v with rates  $1/\lambda(t)$  and  $Q^{\rm f}(t)/\gamma$ , respectively. On the other hand,  $W^{\rm f}(t)$  decreases with the completion of v streams and with requests leaving  $W^{\rm f}(t)$  to enter  $Q^{\rm f}(t)$  or leave the system permanently with rates  $r(t)/\kappa$  and  $1/\alpha$ , respectively. The queue length  $Q^{\rm f}(t)$  increases when requests leave  $Q^{\rm f}(t)$  to enter  $W^{\rm f}(t)$  and decreases when requests leave  $Q^{\rm f}(t)$  to enter  $W^{\rm f}(t)$  with rates  $(1 - \phi)/\alpha$  and  $1/\gamma$ , respectively. The queue lengths  $W^{\rm f}(t)$  and  $Q^{\rm f}(t)$  are used for deriving the delay as discussed in Section II-G.

#### D. User mobility model

The user mobility model discussed below is adopted from [54]. It is assumed that at  $t_1 = 0$ , the user is at a given location. After waiting for a random time  $\Delta t$ , the user moves to a new location with probability  $\rho t^{-\alpha_1\nu/(1+\nu)}$  and otherwise remains in the same location. Thereafter, the process repeats with another random time  $\Delta t$ . The time interval  $\Delta t$  is chosen from a heavy-tailed distribution  $\propto |\Delta t|^{-1-\alpha_1}$ . When the user moves, the distance  $\Delta r$  is chosen from another heavy-tailed distribution  $\propto |\Delta t|^{-1-\alpha_1}$ . The time interval  $\Delta t$  is uniform in  $[0, 2\pi)$ . The mean-square user displacement with respect to their initial position is

$$E[\Delta r^2] \sim \left[ \log \left( \frac{1 - S(t_1)^{1-\nu}}{\nu - 1} \right) \right]^{2/\alpha_2} + \text{constant.} \quad (11)$$

Moreover, the number of users in a particular geographical region is assumed to have a periodic variation [3]. Thus we assume

$$N(t) = \frac{N_{\rm h} - N_{\rm l}}{2} \sin\left(\frac{2\pi t}{T_N}\right) + \frac{N_{\rm h} + N_{\rm l}}{2}.$$
 (12)

#### E. Video popularity model

The request inter-arrival time of a video is measured in real time by measuring the time elapsed from arrival of the previous request from all users at the CPU [36]. Video popularity is modeled by  $\lambda(t)$  assuming direct and word of mouth recommendations  $\alpha_4$  and  $\alpha_5$ , respectively, videos' intrinsic popularity q, and total number of users at time t N(t), resulting in [55]

$$\lambda(t) = \frac{qN(t)\frac{\alpha_4}{\alpha_5 qN(t)}e^{(\alpha_4 + \alpha_5 qN(t))t} - \frac{\alpha_4}{\alpha_5}}{\frac{\alpha_4}{\alpha_5 qN(t)}e^{(\alpha_4 + \alpha_5 qN(t))t} + 1}.$$
 (13)

Equations (4.1)–(4.7) of [56] can be used to derive  $\lambda(t)$  at the CPU from real-time video request inter-arrival time data.

#### F. Channel model

The channel between the UAP and a user consists of line-ofsight and non-line-of-sight links. Assuming elevation of a UAP h, carrier frequency f, speed of light c, path loss exponent  $\alpha_3$ , excessive path loss co-efficient for line-of-sight links  $\delta_1$  and non-line-of-sight links  $\delta_2$ , horizontal distance of the user from the UAP d(t), bandwidth B, environment-dependent constants  $\psi_1$  and  $\psi_2$ , additive white Gaussian noise power  $\sigma^2$ , and UAPto-user transmit power p(t), the probabilities of the line-ofsight and non-line-of-sight links are

$$p_{\rm LoS}(t) = \frac{1}{1 + \psi \exp\left(-\psi_1\left[\frac{180}{\pi}\arctan(\frac{h}{d(t)}) - \psi\right]\right)}$$
(14)

and  $p_{\text{NLoS}}(t) = 1 - p_{\text{LoS}}(t)$ , respectively [57]–[60]. The path loss between a UAP and user is given by

$$P(t) = \left[\frac{4\pi fh}{c}\right]^{\alpha_3} \left(\delta_1 p_{\rm LoS}(t) + \delta_2 p_{\rm NLoS}(t)\right).$$
(15)

The CPU C uses the worst-case interference  $\Xi$  to estimate the download rate r(t). For this purpose, C assumes that a user receives interference from all UAPs except its serving UAP. Moreover, other users are assumed to be present within  $\Delta r$  of a user while the UAPs are positioned to serve these users. Assuming signal power used by UAPs p(t) and path loss between an interfering UAP i and the considered user  $P_i(t)$ ,  $\Xi$  is derived similar to equation (5) of [35] as

$$\Xi = \sum_{i=1}^{\zeta(t)-1} \frac{p(t)}{P_i(t)}.$$
(16)

The path losses  $P_i(t)$  in (16) are derived similarly to P(t) in (15) for real-time positions of UAPs and users. For simplicity of analysis, C is assumed to estimate the same interference for all users. On the other hand, each UAP maximizes the signal-to-interference-plus-noise ratio in real time using the interference-management equations (6)–(28) of [35]. This results in higher streaming rates compared to that estimated by C. The Doppler shift effect due to mobility of the UAPs and users is assumed to be perfectly compensated. As an optimistic estimate of the time-varying download rate, we use

$$r(t) = B \log_2 \left( 1 + \frac{p(t)}{P(t)(\sigma^2 + \Xi)} \right), \tag{17}$$

at C which equals the channel capacity if p(t) and P(t) vary slowly enough.

## *G.* Performance parameter: Delay $\delta^{f}(t)$

Delay is experienced by users when the number of requests for v arriving from users is greater than the number of requests being currently streamed by UAPs. In this paper, the delay is defined in terms of the virtual delay, which is the waiting time experienced by a "test" user that requests v exogenously at any instant t, assuming that the channel conditions do not deteriorate before a UAP starts to stream v to it [51, Sec. 3.3]. In this section, we derive an expression for the delay. We first define the diffusion processes corresponding to the length of the service and retrial queues and thereafter assume that  $\Omega(t)$  is observed in the absence of exogenous arrivals. Delay is defined as an attainment process, i.e., the first instant at which a waiting request begins to stream v. For this purpose we make the following assumptions.

- 1)  $\zeta(t)$  is continuously differentiable at all  $t \ge 0$ .
- 2) r(t) is continuous at all  $t \ge 0$ .
- No new exogenous requests arrive after the arrival of a request from a "test" user at time τ ≥ 0.

Furthermore, we assume that the limit  $\lim_{\eta\to\infty} \sqrt{\eta}(\Omega^{\eta}(0)/\eta - \Omega^{f}(0))$  exists and converges in distribution to  $\Omega^{d}(0) = (W^{d}(0), Q^{d}(0))$  for  $\eta \to \infty$ . The fluid approximation  $\Omega^{f}(t)$  can be refined using the functional central limit theorem [53]. Thus, for t > 0 and the diffusion process  $\Omega^{d}(t)$  [61, Sec. 3.1], [49]

$$\lim_{\eta \to \infty} \sqrt{\eta} \left( \frac{1}{\eta} \mathbf{\Omega}^{\eta}(t) - \mathbf{\Omega}^{\mathrm{f}}(t) \right) \xrightarrow{d} \mathbf{\Omega}^{\mathrm{d}}(t), \tag{18}$$

is a convergence in distribution  $\stackrel{d}{\rightarrow}$  of the stochastic processes in an appropriate functional space [53]. Moreover, if the set of time points  $\{t \ge 0 \mid W^{\rm f}(t) = \zeta(t)\}$  has measure zero for the given queueing system, then  $\Omega^{\rm d}(t)$  is a Gaussian process [49] for  $t \ge 0$ .

The arrival process of v requests to C and the departure process associated with the streaming of v from W(t) are denoted by A(t) and  $\Delta(t)$  with scaled processes  $A^{\eta}(t)$  and  $\Delta^{\eta}(t)$ , respectively, for  $t \ge 0$ . In order to calculate the delay,  $\Omega(t)$  is observed after time  $t = \tau$  for which the queues are denoted by the operator  $\widehat{}$ . Since there are no exogenous future arrivals,  $A^{\eta}(0) = \widehat{W}^{\eta}(0)$ ,  $\Delta^{\eta}(0) = 0$ ,  $A^{\eta}(t) - \Delta^{\eta}(t) = \widehat{W}^{\eta}(t)$ for  $t \ge 0$  [62, Chapter 2], and  $A^{f}(t) = A^{f}(\tau)$  for  $t \ge \tau$ . The fluid limit results of (5) and (6) are written for  $t \ge \tau$  as

$$\lim_{\eta \to \infty} \frac{1}{\eta} \left( \widehat{\mathbf{\Omega}}^{\eta}(t), A^{\eta}(t), \Delta^{\eta}(t) \right) \xrightarrow{d} \left( \widehat{\mathbf{\Omega}}^{\mathrm{f}}(t), A^{\mathrm{f}}(t), \Delta^{\mathrm{f}}(t) \right).$$
(19)

 $\Delta^{\mathrm{f}}(t)$  is a continuously differentiable nondecreasing function of t for  $t \geq 0$ . For  $0 \leq t \leq \tau$ ,  $\widehat{W}^{\mathrm{f}}(t)$  satisfies (9) while for  $t \geq \tau$ , it is assumed that C does not have any new arrivals (exogenous or retrial) after time  $\tau$ . Removing the terms corresponding to the arrivals after time  $\tau$  from (9) results in

$$\frac{\mathrm{d}}{\mathrm{d}t}\widehat{W}^{\mathrm{f}}(t) = -\frac{\min\left(\widehat{W}^{\mathrm{f}}(t),\zeta(t)\right)r(t)}{\kappa} - \frac{\left(\widehat{W}^{\mathrm{f}}(t)-\zeta(t)\right)^{+}}{\alpha}.$$
(20)

The diffusion limits, which are obtained from the convergence in distribution similarly to (18), [61, Sec. 3.1], are

$$\lim_{\eta \to \infty} \sqrt{\eta} \left( \frac{1}{\eta} \widehat{\boldsymbol{\Omega}}^{\eta}(t) - \widehat{\boldsymbol{\Omega}}^{\mathrm{f}}(t), \frac{1}{\eta} A^{\eta}(t) - A^{\mathrm{f}}(t), \frac{1}{\eta} \Delta^{\eta}(t) - \Delta^{\mathrm{f}}(t) \right) \xrightarrow{d} \left( \widehat{\boldsymbol{\Omega}}^{\mathrm{d}}(t), A^{\mathrm{d}}(t), \Delta^{\mathrm{d}}(t) \right).$$
(21)

From (18), (21), and [62, Chapter 2], the diffusion limit  $\widehat{W}^{d}(t)$  can be written as

$$\widehat{W}^{\mathsf{d}}(t) = A^{\mathsf{d}}(t) - \Delta^{\mathsf{d}}(t).$$
(22)

We define the  $\eta$ -scaled process  $E^{\eta}(t)$  as the sum of the processes corresponding to the completion of streaming and ongoing v streams as

$$E^{\eta}(t) = \Delta^{\eta}(t) + \eta \zeta(t), \qquad t \ge 0.$$
(23)

It follows that  $(1/\eta)E^{\eta}(t)$  converges to  $E^{\rm f}(t)$  almost surely for  $\eta \to \infty$  where the convergence is uniform on compact sets of t and  $E^{\rm f}(t) = \Delta^{\rm f}(t) + \zeta(t)$ ,  $t \ge 0$ . Since  $\zeta(t)$  and  $\Delta^{\rm f}(t)$ are continuously differentiable by earlier assumptions,  $E^{\rm f}(t)$ is also continuously differentiable with derivative  $e^{\rm f}(t)$ . It is assumed that  $e^{\rm f}(t)$  is strictly positive and  $\lim_{t\to\infty} E^{\rm f}(t) > A^{\rm f}(\tau)$ . According to previous definitions,  $A^{\eta}(t)$  and  $A^{\rm f}(t)$  are constant for  $t \ge \tau$  assuming that the processes considered above are defined for  $t \ge -T$  with

$$T = \frac{\zeta(0)}{e^{f}(0)}.$$
 (24)

This extension is made by assuming that nothing is happening in  $-T \le t < 0$  (no arrivals or departures) except that the number of users simultaneously streaming v increases linearly from 0 to  $\eta\zeta(0)$ . Therefore, (19) and (21) are retained without modification with all the functions being defined for t > -T.

The processes  $A^{f}(t)$ ,  $E^{f}(t)$ ,  $A^{d}(t)$ , and  $E^{d}(t)$  are continuous and  $E^{f}(-T) = E^{d}(-T) = 0$ , where  $\Delta^{d}(t) = E^{d}(t)$ . Therefore, the processes discussed next are finite with probability 1 for all sufficiently large  $\eta$ . The first attainment processes are defined by the first instant s at which the number of times v is streamed exceeds the number of times it is requested after the arrival of request from the "test" user, i.e., the first instant when this request can start to stream v and are defined for all  $t \ge -T$  as [49]

$$P^{\eta}(t) = \min\left\{s \ge -T : E^{\eta}(s) > A^{\eta}(t)\right\}, \qquad (25)$$

$$P^{f}(t) = \min\left\{s \ge -T : E^{f}(s) > A^{f}(t)\right\}.$$
 (26)

The first attainment processes correspond to the first hitting time in a normal queuing process without a virtual queue [63, Sec. 6.2]. Thus,  $P^{f}(t)$ , is the smallest instant  $\tau \ge t$  such that  $\widehat{W}^{f}(\tau) = \zeta(\tau)$ , i.e., the test user starts streaming v at time  $\tau$ [49]. The delay experienced by the test user is given by the attainment waiting time processes as [51, Sec. 3.3]

$$\delta^{\eta}(t) = P^{\eta}(t) - t, \qquad (27)$$

$$\delta^{\rm f}(t) = P^{\rm f}(t) - t. \tag{28}$$

Moreover,  $\hat{\delta}^{\eta}(\tau)$  denotes the virtual waiting time at  $\tau$ , i.e., the delay experienced by the test user arriving to the service node at time  $\tau$ , until it starts streaming v, assuming that this user does not quit or enter Q(t) while waiting. Then the relation between the virtual waiting time and the attainment waiting time is  $\hat{\delta}^{\eta}(t) = \delta^{\eta}(t)^+$  for all  $t \ge 0$ . Furthermore, the following convergences in distribution follow from (19), (21), [61, Sec. 3.1], [49]

$$\lim_{\eta \to \infty} \left( \frac{1}{\eta} \widehat{\mathbf{\Omega}}^{\eta}(t), \frac{1}{\eta} A^{\eta}(t), \frac{1}{\eta} E^{\eta}(t), \frac{1}{\eta} \delta^{\eta}(t) \right) \\ \stackrel{d}{\to} \left( \widehat{\mathbf{\Omega}}^{\mathrm{f}}(t), A^{\mathrm{f}}(t), E^{\mathrm{f}}(t), \delta^{\mathrm{f}}(t) \right),$$
(29)

$$\lim_{\eta \to \infty} \sqrt{\eta} \left( \frac{1}{\eta} \widehat{\mathbf{\Omega}}^{\eta}(t) - \widehat{\mathbf{\Omega}}^{\mathrm{f}}(t), \frac{1}{\eta} A^{\eta}(t) - A^{\mathrm{f}}(t), \\ \frac{1}{\eta} E^{\eta}(t) - E^{\mathrm{f}}(t), \delta^{\eta}(t) - \delta^{\mathrm{f}}(t) \right) \\ \stackrel{d}{\to} \left( \widehat{\mathbf{\Omega}}^{\mathrm{d}}(t), A^{\mathrm{d}}(t), E^{\mathrm{d}}(t), \delta^{\mathrm{d}}(t) \right),$$
(30)

where

$$\delta^{\mathbf{d}}(t) = \frac{A^{\mathbf{d}}(t) - E^{\mathbf{d}}\left(P^{\mathbf{f}}(t)\right)}{e^{\mathbf{f}}\left(P^{\mathbf{f}}(t)\right)} \tag{31}$$

is derived in [49]. Since the processes  $A^{d}(t)$ ,  $E^{d}(t)$ ,  $\hat{Q}^{d}(t)$ , and  $\delta^{d}(t)$  are continuous, the finite-dimensional distributions converge. In particular, considering the nontrivial case  $P^{f}(\tau) \geq \tau$  (i.e.,  $\widehat{W}^{f}(\tau) \geq \zeta(\tau)$ ), if  $0 \leq t \leq \tau$ , then the set of points  $\{t \mid \widehat{W}^{f}(t) = \zeta(t)\}$  has measure zero, and  $\delta^{\eta}(t)$  converges to  $\delta^{f}(t)$  almost surely for  $\eta \to \infty$ . Thus, from (22), (31), the definition of  $\Delta^{d}(t)$ , and assumption 3 above that results in  $A^{1}(t) = A^{1}(P^{f}(t))$  for all t > -T,

$$\delta^{\mathbf{d}}(t) \xrightarrow{d} \frac{\widehat{W}^{\mathbf{d}}\left(P^{\mathbf{f}}(t)\right)}{e^{\mathbf{f}}\left(P^{\mathbf{f}}(t)\right)}.$$
(32)

Also,  $e^{f}(P^{f}(t)) = \zeta(P^{f}(t))r(P^{f}(t))/\kappa$  when  $P^{f}(t) \ge t$ .

Since a delay is experienced by users when the number of v requests arriving from the request arrival process and waiting to stream v,  $\widehat{W}^{f}(\tau)$  is larger than the requests being streamed  $\zeta(\tau)$ .  $P^{f}(t)$  from (26) can be written as the first instant when the above occurs, i.e.,

$$P^{\rm f}(t) = \min\left\{\tau \ge t \mid \widehat{W}^{\rm f}(\tau) = \zeta(\tau)\right\}.$$
(33)

To compute  $P^{f}(t)$ , (9) is used for  $t < \tau$  and (20) for  $t \ge \tau$ .

A delay is experienced by the users if the UAPs simultaneously stream to  $\zeta(t)$  users only, i.e.,  $\widehat{W}^{f}(t) \geq \zeta(t)$ . In such a scenario, the fluid limit of the delay  $\delta^{f}(t)$  is the difference between the instants when the request from the test user arrived and when a UAP started to stream it,  $P^{f}(t)$ . Therefore, the delay is given by  $\delta^{f}(t) = P^{f}(t) - t$  in analogy with (27).

# III. RESULTS

In this section, we conduct simulations using the queueing model derived in Section II for the scenarios depicted in Table III, and compare them with simulations carried out in OMNET++. Our simulations illustrate that the analytical models are fairly accurate in the asymptotic regime considered in this work. Further, using an RL approach, we also provide guidelines for resource-efficient deployment strategies of UAPs and CPU design.

A maximum of I = 144 UAPs are assumed to serve up to 50000 users in [16]. The period of simulation is limited to 50 s. The variation in the number of video requests arriving as a function of time is illustrated in [65]. The values of parameters for user mobility, channel model, and video popularity model are taken from [64], [57]–[60], and [55], respectively. The UAP-to-user transmit power p(t) is assumed to be a constant 10 mW [58]. The values of the parameters that are common in all simulations are grouped under the scenario S and listed in Table III. The performance results for S are a benchmark in the calculations of the delay. In scenario  $S_1$ ,  $\kappa$  is increased

TABLE III: Values assumed for the parameters [55], [57]–[60], [64]

Sc.	Parameter	Value	Sc.	Parameter Value			
S	$ \begin{matrix} \kappa \\ N_{\rm h} \\ \phi \\ \rho \\ \alpha_1, \alpha_2 \\ \alpha_3 \\ \alpha_4, \alpha_5 \\ \psi_1 \\ \delta_2 \\ B \end{matrix} $	$ \begin{array}{c} 10^{6} \text{ bits} \\ 50000 \\ 0.5 \\ 0.6 \\ 0.55, 0.8 \\ 4 \\ 0.1, 10^{-7} \\ 0.43 \\ 21 \\ 1 \text{ MHz} \\ 0.5, 0.45, 0.05 \end{array} $	$oldsymbol{S}$	$\begin{array}{c} T_N \\ N_1 \\ c \\ h \\ \nu \\ q \\ \psi \\ \delta_1 \\ f_c \\ \sigma^2 \\ \delta_M \end{array}$	$\begin{array}{c} 50\\1\\3\cdot10^8\text{m/s}\\10\text{ m}\\1.21\\0.05\\4.88\\0.01\\2\text{ GHz}\\10^{-13}\\20\text{ s}\\\end{array}$		
	$a_3, a_4, a_5$	0.5, 0.45, 0.05		q	0.5		
$oldsymbol{S}_1$	$\kappa \nearrow$	$10^9$ bits	$oldsymbol{S}_2$	h↗	100 m		
$oldsymbol{S}_3$	$a_3 \nearrow$	0.9, 0.05, 0.05	$oldsymbol{S}_4$	$B \nearrow$	50 MHz		
	$egin{array}{ccc} a_4 &\searrow & \ a_5 & & \ \end{array}$						
$oldsymbol{S}_5$	$T_N \nearrow$	100	$oldsymbol{S}_6$	$N_{\rm h} \searrow$	1000		
$oldsymbol{S}_7$	$\nu \nearrow$	2	$oldsymbol{S}_8$	$I \nearrow$	200		
$oldsymbol{S}_9$	$q \searrow$	$5.10^{-5}$					

to observe the effect if a more bandwidth-intensive video is requested. The height h of UAPs is increased in  $S_2$  to observe the effect of interference while UAPs implement interference management. The effect of a different deployment strategy is illustrated in  $S_3$ . In  $S_4$ , the effect of increased link bandwidth B is explored. The user arrival pattern is changed by increasing the period of variation, decreasing the maximum number of users in a geographical area, and decreasing the users' tendency to explore a new region in  $S_5$ - $S_7$ , respectively. The maximum number of UAPs deployable is increased in  $S_8$ while the effect of a video with low intrinsic popularity is explored in  $S_9$ .

The delay experienced in  $S-S_9$  is reported in Figs. 4(a)-4(j). The analytical model is first solved numerically using MATLAB assuming a periodic deployment of UAPs of the form  $\zeta(t) = \frac{\zeta_{\rm h} - \zeta_{\rm l}}{2} \sin(\frac{2\pi t}{T_{\zeta}}) + \frac{\zeta_{\rm h} + \zeta_{\rm l}}{2}$  with  $T_{\zeta} = 50$ s. The delay obtained from the analytical model proposed in Section II-C is validated using network simulations in OMNET++, assuming that the user base of C is limited to at most 50000 users. In the network simulations, the delay is measured as the difference in the time instants when a user requests vand some UAP starts to stream it. We observe from Fig. 4 that the qualitative nature of the delay obtained from network simulations and the analytical model match well. This validates the fluid approximation models derived in Section II. Thus, our proposed analytical framework is fairly accurate in the asymptotic regime considered in this work, and can aid analysis of the real-time QoS in UANs. However, a UAN should employ resource-efficient strategies to achieve a desired QoS, which requires involvement of the underlying model dynamics. We implement an RL agent that utilizes our fluid approximation model to achieve the above objectives.

The following subsections discuss the effect of varying the parameters listed in Table III on the delay  $\delta^{f}(t)$  experienced by a user while downloading v from the UAN associated



Fig. 3: Convergence of rewards in S.

with C and the UAP deployment strategies implemented. The discussion proceeds by comparing  $\delta^{f}(t)$  observed in  $S_1-S_9$  with that of S.

#### A. Delay

The requests wait in the service or retrial queue at C before being streamed when  $\widehat{W}^{f}(t) > \zeta(t)$ , resulting in a delay in  $S-S_9$ . The delay increases with N(t), k, h,  $N_h$ ,  $\nu$ , and q as observed from the delay plots of  $S-S_2$ ,  $S_6$ ,  $S_7$ , and  $S_9$ , respectively, while it decreases on increasing B and  $T_N$  in  $S_4$  and  $S_5$ , respectively, as observed from Fig. 4. In  $S-S_9$ ,  $\lambda(t)\kappa/\zeta(t)r(t) > 1$  for certain time intervals, resulting in a very high delay in Fig. 4. Therefore, the request queuing delay is high in many scenarios. Furthermore, a variation in  $\zeta(t)$  affects the delay more than N(t). Thus, an efficient UAN design involves jointly optimizing the available resources and the real-time QoS experienced by the users. However,  $\delta^{f}(t)$ is a non-convex function. An optimization problem involving  $\delta^{\rm f}(t), \gamma$ , and  $\zeta(t)$  is difficult to solve. In the next subsection, we employ an RL agent in the CPU to solve the optimization problem and decide UAP deployment that can yield a high QoS.

#### B. Design of UAP deployment strategies

The CPU can change UAP deployment as a function of user mobility and the popularity of v for improving the QoS experienced by users while maximizing resource efficiency. Varying  $\zeta(t)$  and  $\gamma$ , resulting in changes in the UAP deployment and CPU design, respectively, affects resource utilization. The goal is to maximize resource efficiency by minimizing  $\zeta(t)$  and  $\gamma$ while minimizing delay which is also upper-bounded. Thus, the optimization problem is formulated as

$$\begin{aligned}
 P_0: & \min a_1 \zeta(t) + a_2 \gamma \\
 s.t. & 0 \ge \delta^{\rm f}(t) < \delta_{\rm M}, \\
 & \zeta_1 \le \zeta(t) \le \zeta_{\rm h} \\
 & 0 \le \gamma \le \gamma_{\rm M}
 \end{aligned}$$
(34)

for positive constants  $a_1$ ,  $a_2$ . The problem  $\mathbf{P}_0$  is a non-convex problem. We propose a deep Q network (DQN)-based solution method for (34). The state spaces  $\mathcal{S}(t)$ , action spaces  $\mathcal{A}(t)$ , and reward of the proposed DQN  $\mathcal{R}(t)$  are designed by [25] for positive constant  $\beta_1$ ,  $a_3$ ,  $a_4$ ,  $a_5$ ,  $0 \leq \gamma$ , and  $\zeta(t) \geq \zeta_1$ such that

$$\begin{aligned} \mathcal{S}(t) &= \delta^{\rm f}(t) \\ \mathcal{A}(t) &= \gamma, \zeta(t) \\ \mathcal{R}(t) &= \begin{cases} \beta_1 (1 - \zeta(t)/\zeta_{\rm h}), \text{if } \delta^{\rm f}(t) < \delta_{\rm M} \\ -\beta_1 (a_3\zeta(t)/\zeta_{\rm h} - a_4\delta^{\rm f}(t)/\delta_{\rm M} - a_5\gamma/\gamma_{\rm M}), \\ \text{otherwise.} \end{cases} \end{aligned}$$
(35)

The constants  $a_1$  and  $a_2$  refer to the weights given to the number of UAPs deployed and the mean retrial time, respectively, in the optimization problem of (34). The problem  $\mathbf{P}_0$  is translated to derive the reward function  $\mathcal{R}(t)$  of an RL algorithm in (35). In  $\mathcal{R}(t)$ , the constants  $a_3$ ,  $a_4$ , and  $a_5$  can be used to minimize the number of UAPs deployed, the maximum delay, and the mean retrial time, respectively, as done in  $S_3$ . The RL simulations were performed in TensorFlow using Python 3.9. The RL agent is made to explore for 30000 episodes followed by the exploitation phase. The RL agent runs after each episode of duration 0.05 s. It is observed from Fig. 4 and Table IV that the RL agent can jointly minimize the delay and the number of UAPs required in different scenarios. Convergence of the proposed algorithm is tested numerically by running the RL agent for 200 consecutive epochs. Each epoch is assumed to be of 200 s. The average reward obtained by the RL agent in each epoch is plotted in Fig. 3. The average reward obtained in an epoch is observed to increase initially followed by saturation around -62. Thus, we conclude that the proposed algorithm converges in S. The RL agent is also observed to converge in other scenarios for which the average reward is not plotted due to space restrictions. High k and hin  $S_1$  and  $S_2$ , respectively, result in the requirement for high  $\zeta(t)$ . In such scenarios, the RL agent prevents a high delay by adjusting  $\zeta(t)$  and  $\gamma$ , as observed from Figs. 4(b), 4(c), and Table IV, respectively. A stringent policy for reducing the UAPs deployed  $\zeta(t)$ , causes a marginally higher delay in  $S_3$ and Fig. 4(d) compared to S. A high bandwidth availability allows the RL agent to ensure lower resource consumption and delay in  $S_4$  compared to S. A high period of variation in user arrival  $T_N$ , i.e., a less bursty nature of user arrival and less number of users  $N_h$ , results in lower delay and low  $\zeta(t)$  and  $\gamma$ requirements in  $S_5$  and  $S_6$  compared to S. This is illustrated in Figs. 4(f), 4(g) and in table IV. A high tendency of users to come back to their original locations  $\nu$  results in marginally higher delay in Fig. 4(h) for  $S_7$  with higher  $\zeta(t)$  requirement compared to S. Availability of a higher number of deployable UAPs in  $S_8$  compared to S does not cause any changes in delay,  $\gamma$  or  $\zeta(t)$ , since  $\zeta(t) < I$  in **S**. A low video popularity also causes low delay,  $\gamma$  and  $\zeta(t)$  in  $S_9$ . Fig. 4 illustrates the importance of employing the RL agent in maintaining a good QoS in different network scenarios. Fig. 4 also compares the delay achieved with that in [18], [66]. Delay has been derived as a closed-form expression of the time-averaged parameters in [18] and is observed to be higher compared to that achieved by our RL agent from Fig. 4. In fact, the delay achieved by this scheme in scenarios  $S_1$  and  $S_2$  are in the order of 100s of seconds and could not fit into the plots of Figs. 4(b) and 4(c). On the other hand, the scenario of [66] is a system



Fig. 4: Evolution of delay. Dashed black line:  $\delta^{f}(t)$  from MATLAB. Dotted black line with marker: Delay from network simulations. Black line: Delay from RL. Dashed red line with marker: Queueing delay of [66]. Red line: Queueing delay of [18].

TABLE IV: Values obtained from RL for  $\gamma$  and  $\zeta(t)$ 

	$\mid S$	$oldsymbol{S}_1$	$ S_2 $	$oldsymbol{S}_3$	$oldsymbol{S}_4$	$ S_5 $	$oldsymbol{S}_6$	$old S_7$	$oldsymbol{S}_8$	$old S_9$
$\zeta(t)$	93	94	127	98	59	97	78	112	92	62
$\gamma$ (s)	10	10	10	2.5	5	5	2.5	1.42	5	5

with time-averaged Markovian arrival and service processes with finite servers and calling population. Thus, according to the analysis of [66], the queueing delay experienced by users before starting to stream the video is significantly lower and has an almost constant value. However, the time-averaged delay obtained by [18], [66] does not provide a true picture of the real-time QoS experienced by users, as observed from the delay obtained from MATLAB and network simulations in Fig. 4. This affects the UAP deployment. In contrast, our work captures the real-time queuing delay, i.e., QoS, and provides deployment strategies to ensure a desired QoS with resource efficiency.

# **IV. CONCLUSIONS**

Present-day networks often face a challenge to deliver data-intensive multimedia video to highly mobile users with high QoS. Although UANs can offer high QoS, they can result in high deployment cost and impose logistic challenges due to the limitations imposed by UAPs. The correct UAN parameters must be varied to enable high QoS and costefficient deployment of UAPs simultaneously. This highlights the importance of the analytical model derived in this paper.

#### REFERENCES

 Cisco, "Cisco Annual Report," Tech. Rep., 2020. [Online]. Available: https://www.cisco.com/c/en/us/solutions/collateral/executiveperspectives/annual-internet-report/white-paper-c11-741490.pdf, visited on 10/11/2021

- [2] C. Zhan, H. Hu, X. Sui, Z. Liu, J. Wang, and H. Wang, "Joint resource allocation and 3D aerial trajectory design for video streaming in UAV communication systems," *IEEE Transactions on Circuits and Systems* for Video Technology, vol. 31, no. 8, pp. 3227–3241, 2021.
- [3] J. Lu, S. Wan, X. Chen, Z. Chen, P. Fan, and K. B. Letaief, "Beyond empirical models: Pattern formation driven placement of UAV base stations," *IEEE Transactions on Wireless Communications*, vol. 17, no. 6, pp. 3641–3655, 2018.
- [4] M. Banagar and H. S. Dhillon, "Performance characterization of canonical mobility models in drone cellular networks," *IEEE Transactions on Wireless Communications*, vol. 19, no. 7, pp. 4994–5009, 2019.
- [5] J. Ji, K. Zhu, D. Niyato, and Ran Wang, "Joint cache placement, flight trajectory, and transmission power optimization for multi-UAV assisted wireless networks," *IEEE Transactions on Wireless Communications*, vol. 19, no. 8, pp. 5389–5403, 2020.
- [6] Z. Becvar, P. Mach, J. Plachy, and M. F. P. De Tudela, "Positioning of flying base stations to optimize throughput and energy consumption of mobile devices," in *IEEE Vehicular Technology Conference*, 2019.
- [7] Z. Becvar, M. Vondra, P. Mach, J. Plachy, and D. Gesbert, "Performance of mobile networks with UAVs," *European Wireless Conference*, pp. 261–267, 2017.
- [8] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Optimal transport theory for power-efficient deployment of unmanned aerial vehicles," *IEEE International Conference on Communications*, 2016.
- [9] N. Zhao, R. Yu, L. Fan, Y. Chen, J. Tang, A. Nallanathan, and V. C. M. Leung, "Caching unmanned aerial vehicle-enabled small-cell networks," *IEEE Vehichular Technology Magazine*, vol. 14, pp. 71–79, 2019.
- [10] E. Kalantari, H. Yanikomeroglu, and A. Yongacoglu, "Wireless networks with cache-enabled and backhaul-limited aerial base stations," vol. 19, no. 11, pp. 7363–7376, 2020.
- [11] L. Vigneri, T. Spyropoulos, and C. Barakat, "Low cost video streaming through mobile edge caching: Modelling and optimization," *IEEE Transactions on Mobile Computing*, vol. 18, no. 6, pp. 1302–1315, 2019.
- [12] M. Chen, M. Mozaffari, W. Saad, C. Yin, M. Debbah, and C. S. Hong, "Caching in the sky: Proactive deployment of cache-enabled unmanned aerial vehicles for optimized quality-of-experience," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 5, pp. 1046–1061, 2017.
- [13] C. Di Franco and G. Buttazzo, "Energy-aware coverage path planning of UAVs," *IEEE International Conference on Autonomous Robot Systems* and Competitions, pp. 111–117, 2015.
- [14] J. Ji, K. Zhu, C. Yi, and D. Niyato, "Energy consumption minimization in UAV-assisted mobile-edge computing systems: Joint resource allocation and trajectory design," *IEEE Internet of Things Journal*, vol. 8, no. 10, pp. 8570–8584, 2021.
- [15] M. Alzenad, A. El-Keyi, and H. Yanikomeroglu, "3-D placement of an unmanned aerial vehicle base station for maximum coverage of users with different QoS requirements," *IEEE Wireless Communications Letters*, vol. 7, no. 1, pp. 38–41, 2018.

- [16] A. Fotouhi, M. Ding, and M. Hassan, "Flying drone base stations for macro hotspots," *IEEE Access*, vol. 6, pp. 19530–19539, 2018.
- [17] S. H. Cheng, Y. T. Shih, and K. C. Chang, "Proactive deployment of cache-enabled aerial base stations for optimized energy-delay cost," in *IEEE International Symposium on Personal, Indoor and Mobile Radio Communications, (PIMRC)*, 2022, pp. 493–498.
- [18] S. Kumar, S. Suman, and S. De, "Backhaul and delay-aware placement of UAV-enabled base station," *IEEE Conference on Computer Commu*nications Workshops, pp. 634–639, 2018.
- [19] S. Anokye, D. Ayepah-Mensah, A. M. Seid, G. O. Boateng, and G. Sun, "Deep reinforcement learning-based mobility-aware UAV content caching and placement in mobile edge networks," *IEEE Systems Journal*, vol. 16, no. 1, pp. 275–286, 2022.
- [20] K. Shanmugam, N. Golrezaei, A. G. Dimakis, A. F. Molisch, and G. Caire, "FemtoCaching: Wireless content delivery through distributed caching helpers," *IEEE Transactions on Information Theory*, vol. 59, no. 12, pp. 8402–8413, 2013.
- [21] D. H. Tran, S. Chatzinotas, and B. Ottersten, "Satellite- and cacheassisted UAV: A joint cache placement, resource allocation, and trajectory optimization for 6G aerial networks," *IEEE Open Journal of Vehicular Technology*, vol. 3, pp. 40–54, 2022.
- [22] J. Luo, J. Song, F. C. Zheng, L. Gao, and T. Wang, "User-centric UAV deployment and content placement in cache-enabled multi-UAV networks," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 5, pp. 5656–5660, 2022.
- [23] A. Masood, T. V. Nguyen, T. P. Truong, and S. Cho, "Content caching in HAP-assisted multi-UAV networks using hierarchical federated learning," *International Conference on ICT Convergence*, pp. 1160–1162, 2021.
- [24] L. Zhao, W. Sun, Y. Shi, and J. Liu, "Optimal placement of cloudlets for access delay minimization in SDN-based internet of things networks," *IEEE Internet of Things Journal*, vol. 5, no. 2, pp. 1334–1344, 2018.
- [25] L. Liu, K. Xiong, Y. Lu, P. Fan, and K. B. Letaief, "Age-constrained energy minimization in UAV-assisted wireless powered sensor networks: A DQN-based approach," *IEEE Conference on Computer Communications Workshops*.
- [26] B. Liu, Q. Xie, and E. Modiano, "RL-QN: A reinforcement learning framework for optimal control of queueing systems," in *Annual Allerton Conference on Communication, Control, and Computing*, Allerton, 2019.
- [27] T. Zhang, Z. Wang, Y. Liu, W. Xu, and A. Nallanathan, "Joint resource, deployment, and caching optimization for AR applications in dynamic UAV NOMA networks," *IEEE Transactions on Wireless Communications*, vol. 21, no. 5, pp. 3409–3422, 2022.
- [28] B. Zeng, C. Zhan, Y. Yang, and J. Liao, "Access delay minimization for scalable videos in cache-enabled multi-UAV networks," in *Proceedings* of the IEEE Global Communications Conference, GLOBECOM 2022. IEEE, 2022, pp. 389–394.
- [29] T. Zhang, Y. Wang, Y. Liu, W. Xu, and A. Nallanathan, "Cache-enabling UAV communications: Network deployment and resource allocation," *IEEE Transactions on Wireless Communications*, vol. 19, no. 11, pp. 7470–7483, 2020.
- [30] J. Wu, Y. Zhou, D. M. Chiu, and Z. Zhu, "Modeling dynamics of online video popularity," *IEEE Transactions on Multimedia*, vol. 18, no. 9, pp. 1882–1895, 2016.
- [31] E. Kalantari, H. Yanikomeroglu, and A. Yongacoglu, "Wireless networks with cache-enabled and backhaul-limited aerial base stations," *IEEE Transactions on Wireless Communications*, vol. 19, no. 11, pp. 7363– 7376, 2020.
- [32] E. Björnson and L. Sanguinetti, "Scalable cell-free massive MIMO systems," *IEEE Transactions on Communications*, vol. 68, no. 7, pp. 4247–4261, 2020.
- [33] H. V. Abeywickrama, Y. He, E. Dutkiewicz, B. A. Jayawickrama, and M. Mueck, "A reinforcement learning approach for fair user coverage using UAV mounted base stations under energy constraints," *IEEE Open Journal of Vehicular Technology*, vol. 1, no. January, pp. 67–81, 2020.
- [34] M. Ding and D. L. Perez, "Please lower small cell antenna heights in 5G," 2016 IEEE Global Communications Conference, GLOBECOM 2016 - Proceedings, pp. 2–7, 2016.
- [35] U. Challita, W. Saad, and C. Bettstetter, "Interference management for cellular-connected UAVs: A deep reinforcement learning approach," *IEEE Transactions on Wireless Communications*, vol. 18, no. 4, pp. 2125–2140, 2019.
- [36] T. Wang, C. Jayasundara, M. Zukerman, A. Nirmalathas, E. Wong, C. Ranaweera, C. Xing, and B. Moran, "Estimating video popularity from past request arrival times in a VoD system," *IEEE Access*, vol. 8, pp. 19934–19947, 2020.

- [37] N. Wang, G. Shen, S. K. Bose, and W. Shao, "Zone-based cooperative content caching and delivery for radio access network with mobile edge computing," *IEEE Access*, vol. 7, pp. 4031–4044, 2019.
- [38] R. I. Ansari, N. Ashraf, and C. Politis, "An energy-aware distributed open market model for UAV-assisted communications," *IEEE Vehicular Technology Conference*, vol. 2020-May, 2020.
- [39] Y. K. Tun, Y. M. Park, T. H. T. Le, Z. Han, and C. S. Hong, "A business model for resource sharing in cell-free UAVs-assisted wireless networks," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 8, pp. 8839–8852, 2022.
- [40] N. Wang, G. Shen, S. K. Bose, and W. Shao, "Zone-based cooperative content caching and delivery for radio access network with mobile edge computing," *IEEE Access*, vol. 7, pp. 4031–4044, 2019.
- [41] S. Kourtis and R. Tafazolli, "Evaluation of handover related statistics and the applicability of mobility modelling in their prediction," *IEEE International Symposium on Personal, Indoor and Mobile Radio Communications, (PIMRC)*, pp. 665–670, 2000.
- [42] B. Lin, W. Wang, J. Guo, and Z. Fei, "Outage Performance for UAV Communications under Imperfect Beam Alignment: A Stochastic Geometry Approach," *International Conference on Communication Technology Proceedings, ICCT*, pp. 632–637, 2021.
- [43] F. A. Cruz-Pérez, A. Seguín-Jiménez, and L. Ortigoza-Guerrero, "Effects of handoff margins and shadowing on the residence time in cellular systems with link adaptation," *IEEE Vehicular Technology Conference*, 2004.
- [44] P. K. Sharma and D. I. Kim, "Coverage probability of 3D mobile UAV networks," *IEEE Wireless Communications Letters*, vol. 8, no. 1, pp. 2018–2021, 2018.
- [45] L. Libman and A. Orda, "Optimal retrial and timeout strategies for accessing network resources," *IEEE/ACM Transactions on Networking*, vol. 10, no. 4, pp. 551–564, 2002.
- [46] C. M. Chen, C. W. Lin, and Y. C. Chen, "Cross-layer packet retry limit adaptation for video transport over wireless LANs," *IEEE Transactions* on Circuits and Systems for Video Technology, vol. 20, no. 11, pp. 1448– 1461, 2010.
- [47] J. Goseling and O. Simeone, "Soft-TTL: Time-varying fractional caching," *IEEE Networking Letters*, vol. 1, no. 1, pp. 18–21, 2019.
- [48] M. A. Maddah-Ali and U. Niesen, "Decentralized coded caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Transactions on Networking*, vol. 23, no. 4, pp. 1029–1040, 2015.
- [49] A. Mandelbaum, W. A. Massey, M. I. Reiman, A. Stolyar, and B. Rider, "Queue lengths and waiting times for multiserver queues with abandonment and retrials," *Telecommunication Systems*, vol. 21:2-4, pp. 149– 171, 2002.
- [50] L. C. Drazek, "Intensity estimation for Poisson processes," Dissertation, Master of Science in Statistics, School of Mathematics, The University of Leeds, 2013. [Online]. Available: http://129.11.36.6/ voss/projects/2012-Poisson/Drazek.pdf, visited on 20/10/2022
- [51] M. El-Taha and S. Stidham, Sample-Path Analysis of Queueing Systems, 1st ed. Springer US, 1999.
- [52] J. G. Dai and J. M. Harrison, Processing Networks Fluid Models and Stability. Cambridge University Press, 2020.
- [53] A. Mandelbaum, W. A. Massey, and M. I. Reiman, "Strong approximations for Markovian service networks," *Queueing Systems*, vol. 30, pp. 149–201, 1998.
- [54] C. Song, T. Koren, P. Wang, and A. L. Barabási, "Modelling the scaling properties of human mobility," *Nature Physics*, vol. 6, no. 10, pp. 818– 823, 2010.
- [55] J. Wu, Y. Zhou, D. M. Chiu, and Z. Zhu, "Modeling dynamics of online video popularity," *IEEE Transactions on Multimedia*, vol. 18, no. 9, pp. 1882–1895, 2016.
- [56] J. Bigot, S. Gadat, T. Klein, and C. Marteau, "Intensity estimation of non-homogeneous Poisson processes from shifted trajectories," *Electronic Journal of Statistics*, vol. 7, no. 1, pp. 881–931, 2013.
- [57] S. Ahmed, M. Z. Chowdhury, and Y. M. Jang, "Energy-efficient UAVto-user scheduling to maximize throughput in wireless networks," *IEEE Access*, vol. 8, pp. 21215–21225, 2020.
- [58] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Mobile unmanned aerial vehicles (UAVs) for energy-efficient internet of things communications," *IEEE Transactions on Wireless Communications*, vol. 16, no. 11, pp. 7574–7589, 2017.
- [59] A. Bhowmick, S. D. Roy, and S. Kundu, "Throughput maximization of a UAV assisted CR network with NOMA-based communication and energy-harvesting," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 1, pp. 362–374, 2022.

- [60] K. M. Liao, G. Y. Chen, Y. J. Chen, and Y. F. Chen, "End-to-end delay analysis in mmWave UAV-assisted wireless caching networks," in *IEEE Wireless Communications and Networking Conference Workshops*, 2020.
- [61] G. Weiss, "Lecture 3: Queueing Networks and their Fluid and Diffusion." [Online]. Available: https://web.stanford.edu/class/msande324/handouts/Lecture3.pdf, visited on 10/11/2021
- [62] L. Kleinrock, Queueing Systems. John Wiley & Sons, 1975.
- [63] R. G. Gallager, Stochastic Processing: Theory For Applications. Cambridge University Press, 2013.
- [64] C. Song, Z. Qu, N. Blumm, and A. L. Barabási, "Limits of predictability in human mobility," *Science*, vol. 327, no. 5968, pp. 1018–1021, 2010.
- [65] J. Liu, H. Yan, Y. Li, D. Karamshuk, N. Sastry, D. Wu, and D. Jin, "Discovering and understanding geographical video viewing patterns in urban neighborhoods," *IEEE Transactions on Big Data*, vol. 7, no. 5, pp. 873–884, 2021.
- [66] Q. Liu, T. Xia, L. Cheng, M. Van Eijk, T. Ozcelebi, and Y. Mao, "Deep reinforcement learning for load-balancing aware network control in IoT edge systems," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 6, pp. 1491–1502, 2022.