

# Statistical QoS Provisionings for Wireless Unicast/Multicast of Multi-Layer Video Streams

Qinghe Du, *Student Member, IEEE*, and Xi Zhang, *Senior Member, IEEE*

**Abstract**—Due to the time-varying wireless channels, deterministic quality of service (QoS) is usually difficult to guarantee for real-time multi-layer video transmissions in wireless networks. Consequently, statistical QoS guarantees have become an important alternative in supporting real-time video transmissions. In this paper, we propose an efficient framework to model the statistical delay QoS guarantees, in terms of QoS exponent, effective bandwidth/capacity, and delay-bound violation probability, for multi-layer video transmissions over wireless fading channels. In particular, a separate queue is maintained for each video layer, and the same delay bound and corresponding violation probability threshold are set up for all layers. Applying the effective bandwidth/capacity analyses on the incoming video stream, we obtain a set of QoS exponents for all video layers to effectively characterize this delay QoS requirement. We then develop a set of optimal adaptive transmission schemes to minimize the resource consumption while satisfying the diverse QoS requirements under various scenarios, including video unicast/multicast with and/or without loss tolerance. Simulation results are also presented to demonstrate the impact of statistical QoS provisionings on resource allocations of our proposed adaptive transmission schemes.

**Index Terms**—Mobile multicast, wireless networks, rate control, layered video streaming, statistical QoS guarantees.

## I. INTRODUCTION

RECENTLY, supporting real-time video services with diverse QoS constraints has become one of the essential requirements for wireless communications networks. Consequently, how to efficiently guarantee QoS for video transmission attracts more and more research attention [1]-[14]. However, the unstable wireless environments and the popular layer-structured video signals [2]-[4] impose a great deal of challenges in delay QoS provisionings. Due to the highly-varying wireless channels, the *deterministic* delay QoS requirements are usually hard to guarantee. As a result, *statistical* delay QoS guarantees [9]-[14], in terms of effective bandwidth/capacity and queue-length-bound/delay-bound violation probabilities, have been proposed and demonstrated as the powerful way to characterize delay QoS provisionings for wireless traffics. While many related existing research works mainly focused on the scenarios with single-layer streams [11]-[14], the modern video coding techniques usually

generate layer-structured traffics [2]-[4]. Unfortunately, how to design efficient schemes to support statistical delay QoS for layered video traffics over wireless networks has been neither well understood nor thoroughly studied.

In video transmissions, video source is usually encoded into a number of data layers [2]-[4] in the application protocol layer. By applying layered video coding, the receivers under poorer channel conditions can get only lower video quality, while those under better channel conditions can achieve higher video quality. Although the layered coding techniques are efficient in handling diverse channel conditions, they also raise new challenges for statistical delay QoS guarantees, which are not encountered in single-layer video transmissions. First, we need to keep the synchronous transmissions across different video layers, implying the same delay-bound violation probability for all layers. Second, for multi-layer video stream, it is a natural requirement that different video layers can tolerate different loss levels. Therefore, the scheduling and resource allocation need to be aware of the diverse loss constraints. Third, how to minimize the consumption of scarce wireless-resources while satisfying the specified delay QoS requirements is a widely cited open problem.

Besides the general challenges in statistical delay QoS guarantees for the unicast transmission of layered video, multicasting layered video over wireless networks further complicates the problem significantly due to the heterogeneous channel qualities across multicast receivers at each time instant. Unlike in the wireless multicast, there are relatively more research results for the multicast over wireline networks. A number of multicast protocols were proposed over wireline networks. The authors of [3] developed the efficient receiver-driven layered multicast over the Internet, where the video source is encoded to a hierarchical signal with different layers. Each layer corresponds to a multicast group and multicast receivers can join/leave the group based on their bandwidths. In [5], we proposed a novel flow control scheme for multicast services over the asynchronous transfer mode (ATM) networks. The kernel parts of this scheme include the optimal second-order rate control algorithm and the feedback soft-synchronization protocol [6], [7], which can achieve *scalable* and *adaptive* multicast flow control over bandwidth and buffer occupancies and utilizations. The above designs are shown to be efficient in the wireline networks. However, the multicast strategies in wireline networks cannot be directly applied into wireless networks. This is because highly and rapidly time-varying wireless-channels qualities result in unstable bandwidths and thus unsatisfied loss and delay QoS. For wireless video multicast, at the multicast sender we need to design the transmission

Manuscript received 1 March 2009; revised 1 November 2009. The research reported in this paper was supported in part by the U.S. National Science Foundation CAREER Award under Grant ECS-0348694. As this paper was co-authored by a guest editor of this issue, the review of this manuscript was handled by senior editor Prof. Larry Milstein.

The authors are with the Networking and Information Systems Laboratory, Department of Electrical and Computer Engineering, Texas A&M University, College Station, TX 77843 USA (e-mails: duqinghe@tamu.edu; xizhang@ece.tamu.edu).

Digital Object Identifier 10.1109/JSAC.2010.100413.

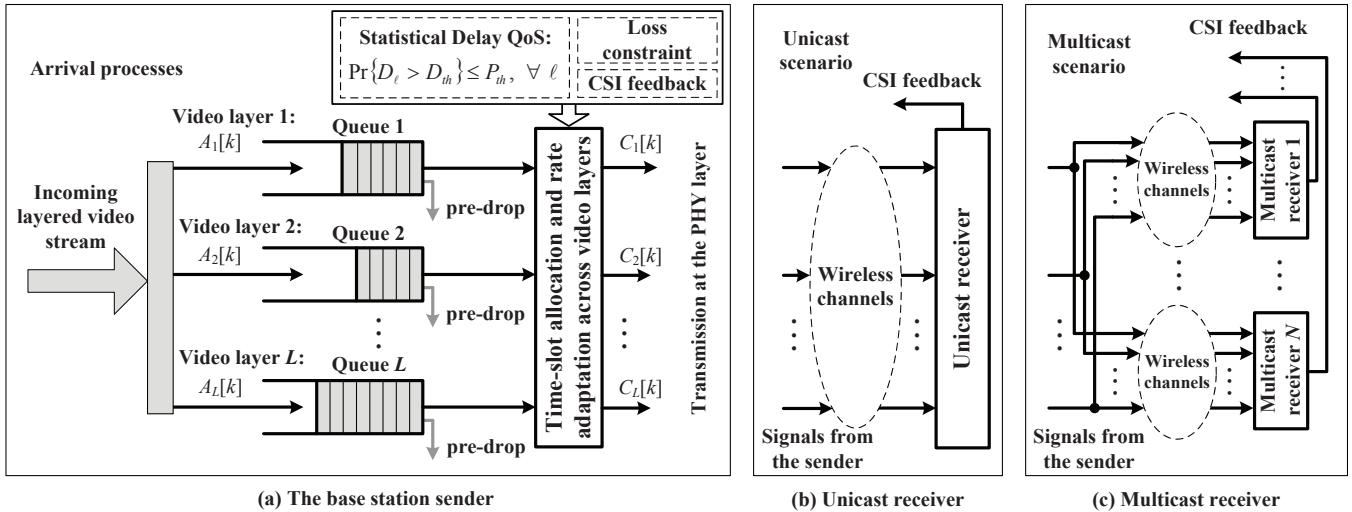


Fig. 1. The system modeling framework for layered-video transmission over wireless networks; (a) The layered-video arrival stream and the sender's processing. (b) Unicast scenario. (c) Multicast scenario.

scheme to control the loss and/or delay performance for all multicast receivers at each video layer based on their instantaneous channel qualities.

In [8] and [14], we applied the effective capacity theory to propose and evaluate rate-adaptation schemes for statistical delay QoS guarantees in mobile multicast. However, the analyses only focused on single-layer stream. It remains one of the major challenges to extend the statistical QoS theory into multi-layer video multicast in developing QoS-driven transmission strategies. In [4], the authors proposed a cross-layer architecture for adaptive video multicast over multirate wireless LANs. In particular, two-layer video signals are considered, which include the base layer (more important) and enhancement layer (less important). The authors derived the transmission rate for the base layer according to the worst-case signal-to-noise ratio (SNR) among all receivers, while dynamically regulating the transmission rate for the enhancement layer based on the best-case SNR to benefit receivers with better channel qualities. However, under this strategy, the loss rate of the enhancement layer will vary significantly with the statistical characteristics of wireless channels, and thus is hard to control.

To overcome the above problems, in this paper we propose an efficient framework to model the statistical delay QoS guarantees, in terms of QoS exponent, effective bandwidth/capacity, and delay-bound violation probability, for multi-layer video transmission over wireless networks. In particular, a separate queue is maintained for each video layer, and the same delay bound and the corresponding violation probability are set up for all video layers. We then develop a set of optimal adaptive transmission schemes to minimize the resource consumption while satisfying the diverse QoS requirements under various scenarios, including video unicast/multicast with and/or without loss tolerance.

The rest of this paper is organized as follows. Section II describes the system model. Section III proposes the framework of statistical delay QoS guarantees for multi-layer video unicast and multicast. Section IV presents the design procedures

for multi-layer video by applying effective bandwidth/capacity theory. Sections V and VI derive the optimal adaptive transmission schemes for video unicast and multicast, respectively. Section VII presents the simulation evaluations. The paper concludes with Section VIII.

## II. THE SYSTEM MODEL

We consider the unicast/multicast system model for multi-layer video distribution in wireless networks, as shown in Fig. 1. Specifically, the base station sender is responsible for transmitting a multi-layer video stream to a single receiver (unicast) or multiple receivers (multicast) over broadcast fading channels. The video stream generated by upper protocol layers (e.g., application layer) consists of  $L$  video layers, each having the specific QoS requirements. The  $L$ -layer video stream will be injected to the physical (PHY) layer. Then, as depicted in Fig. 1(a), we aim at developing strategies to efficiently allocate limited wireless resources for multi-layer video transmission while satisfying the specified QoS requirements for each video layer.

At the PHY layer, the sender uses a constant transmit power with the signal bandwidth equal to  $B$  Hz. The wireless broadcast channels are assumed to be flat fading. Then, we can use an SNR vector  $\gamma \triangleq (\gamma_1, \gamma_2, \dots, \gamma_N)$  to characterize the channel state information (CSI) of receivers, where  $N$  denotes the number of unicast/multicast receivers,  $\gamma_n$  is the received SNR of the  $n$ th receiver for  $n = 1, 2, \dots, N$ , and  $\{\gamma_n\}_{n=1}^N$  are independent and identically distributed (i.i.d.) for the cases of  $N > 1$ . When  $N$  is equal to 1, the scenario reduces to video unicast,<sup>1</sup> as illustrated in Fig. 1(b); while  $N$  is larger than 1, we get the multicast scenario depicted in Fig. 1(c). The CSI  $\gamma$  is modeled as an ergodic and stationary block-fading process, where  $\gamma$  does not change within a time-frame with the fixed length  $T$ , but varies independently from frame to frame. Moreover,  $\gamma$  follows Rayleigh fading model, which is one of the most generally used models to characterize wireless fading channels. In addition, we assume that  $\gamma$  can

<sup>1</sup>When  $N = 1$ , we write SNR as  $\gamma$  instead of  $\gamma_1$  to simplify notation.

be perfectly estimated by the receivers and reliably fed back to the sender without delay through the dedicated feedback control channels.

### III. MODELING FRAMEWORK FOR WIRELESS UNICAST/MULTICAST OF MULTI-LAYER-VIDEO

We propose the following framework for transmitting multi-layer video over fading channels by integrating the adaptive resource allocations, statistical QoS guarantees, and loss constraints.

#### A. Multi-Queue Model for Multi-Layer Video Arrival Processes

The modern video coding techniques [2] usually encode the video source into a number of video layers with different relevance and importance. The most important layer is called based layer and the other layers are called enhancement layers. Because of the diverse importance, different strategies need to be proposed for the corresponding video layers in the PHY-layer transmission, depending on the specified QoS requirements (to be detailed in Sections III-B and III-D). Then, to achieve the efficient video transmission, the sender manages a separate queue for each video layer. As shown in Fig. 1(a), the data arrival rate of the  $\ell$ th video layer is characterized by a discrete-time process, denoted by  $A_\ell[k]$  (nats/frame), where  $\ell = 1, 2, \dots, L$ , and  $[k]$ ,  $k = 1, 2, \dots$ , is the index of time frames; the service rate process (departure process) of the  $\ell$ th layer is denoted by  $C_\ell[k]$  (nats/frame). Moreover, we determine  $C_\ell[k]$  based on CSI, total available wireless resources, and QoS constraints.

#### B. Statistical Delay QoS Guarantees for Video Transmissions

For video transmissions, delay is one of the most important QoS metrics. However, due to the highly varying wireless channels, usually the hard delay bound cannot be guaranteed. Therefore, the statistical metric, namely *delay-bound violation probability* [9]-[12], has been widely applied in QoS evaluations for real-time services. In our framework, we also use the delay-bound violation probability to statistically characterize the delay QoS provisionings for each video layer. In particular, a queueing delay bound, denoted by  $D_{\text{th}}$ , is specified. Accordingly, over all video layers, the delay-bound violation probability cannot exceed the threshold denoted by  $P_{\text{th}}$ :

$$\Pr\{D_\ell > D_{\text{th}}\} \leq P_{\text{th}}, \quad P_{\text{th}} \in (0, 1), \quad (1)$$

for all  $\ell \in \{1, 2, \dots, L\}$ , where  $D_\ell$  denotes the queueing delay at the  $\ell$ th video layer. The delay-bound violation probability in Eq. (1) is evaluated over the entire transmission process, which is assumed to be long enough. Note that  $D_{\text{th}}$  and  $P_{\text{th}}$  are application-dependent parameters. Moreover, the pair of  $D_{\text{th}}$  and  $P_{\text{th}}$  is set to be the same for all video layers, because the synchronous transmission across different video layers is usually required. Also note that the delay in wireless video transmissions may result from multiple factors such as transmissions, queueing, and decoding. In this paper, we mainly focus on queueing delay, which reflects the capability of the wireless channel (transmission bottleneck) in supporting video distributions.

#### C. Adaptive Resource Allocation and Transmissions

To efficiently use the limited wireless resources for video unicast/multicast, we employ the adaptive transmission strategy (based on the CSI), consisting of three folds: transmission rate adaptation, dynamic time-slot allocation, and adaptive pre-drop queue management strategy, as detailed below.

1) *Time Slot Allocation for Video Layers*: Each time frame is divided into  $L$  time slots, the lengths of which are denoted by  $\{T_\ell[k]\}_{\ell=1}^L$ , where  $0 \leq T_\ell[k] \leq T$  and  $\sum_{\ell=1}^L T_\ell[k] \leq T$ . The time slot with length  $T_\ell[k]$  is used for transmitting data of the  $\ell$ th video layer. For convenience of presentation, we further define time proportion  $t_\ell[k] \triangleq T_\ell[k]/T$ , and thus we have  $\sum_{\ell=1}^L t_\ell[k] \leq 1$ . Notice that our target is to minimize the wireless-resource consumption while satisfying the QoS requirements imposed by video qualities. Thus,  $\sum_{\ell=1}^L t_\ell[k]$  may be smaller than 1 for some  $\gamma$ .

2) *Rate Adaptation of Unicast/Multicast*: We denote the total amount of transmitted data at the  $\ell$ th video layer in the  $k$ th time frame by  $\mathcal{R}_\ell[k]$  (with the unit nats/frame). Moreover, we use the normalized transmission rate, denoted by  $R_\ell[k]$  (nats/s/Hz), to characterize the transmission rate adaptation, where  $R_\ell[k] \triangleq \mathcal{R}_\ell[k]/(BT)$ . We assume that capacity-achieving codes are used for transmission at the PHY layer. Accordingly, for unicast, the normalized transmission rate of the  $\ell$ th video layer is set equal to the Shannon capacity under the current SNR  $\gamma$ :

$$R_\ell[k] = \log(1 + \gamma) \quad (\text{nats/s/Hz}). \quad (2)$$

Clearly,  $R_\ell[k]$  does not vary with  $\ell$ , and thus we only focus on time-slot allocation for unicast.

For the multicast case, the rate adaptation becomes more complicated. In particular, the time slot for video layer  $\ell$  is further partitioned into  $N$  sub-slots. The length of the  $n$ th sub-slot, denoted by  $T_{\ell,n}[k]$ , is equal to  $T_\ell[k]t_{\ell,n}[k]$ , where  $0 \leq t_{\ell,n}[k] \leq 1$  and  $\sum_{n=1}^N t_{\ell,n}[k] = 1$ . Within the  $n$ th sub-slot, the transmission rate is set equal to the Shannon capacity under SNR  $\gamma_n$ , and thus the data transmitted in this sub-slot can be correctly decoded only by receivers with SNR higher than or equal to  $\gamma_n$ . Then, the normalized transmission rate  $R_\ell[k]$  for the  $\ell$ th video layer becomes

$$\begin{aligned} R_\ell[k] &= \sum_{n=1}^N t_{\ell,n}[k] R_{\ell,n}[k] \\ &= \sum_{n=1}^N t_{\ell,n}[k] \log(1 + \gamma_n) \quad (\text{nats/s/Hz}), \end{aligned} \quad (3)$$

where  $R_{\ell,n}[k] \triangleq \log(1 + \gamma_n)$  is the normalized transmission rate for the  $n$ th sub-slot of the  $\ell$ th video layer. As a result, we need to not only adjust  $t_\ell[k]$ 's for each layer, but also regulate  $t_{\ell,n}[k]$ 's within every time slot.

Unlike the wireline multicast networks, in this paper we focus on the layered video transmissions over wireless networks, which has a single-hop cellular network structure. Due to the broadcast nature of wireless channels, the sender only needs to transmit a *single copy* of data and *all* multicast receivers can hear the transmitted signal for each video layer. Under this model, our scheme employs the *sender-oriented* multicast approach because the sender needs to dynamically adjust the

transmissions rate in controlling the loss rate (to be detailed in Section III-D) and guaranteeing the delay-QoS requirements (see Section III-B) across different multicast receivers.

3) *Pre-drop Strategy*: In [14], for multicasting single-layer-data we developed the pre-drop strategy to gain a more robust queueing behavior. In this paper, we further extend the pre-drop strategy to multi-layer video transmission. Specifically, based on the CSI, in each time frame the sender can drop some data (see Fig. 1(a)) from the head of each queue, but treat them as if they were transmitted. We denote the amount of dropped data in  $\ell$ th video layer by  $Z_\ell[k]$  (nats/frame) and define the normalized drop rate, denoted by  $z_\ell[k]$ , as  $z_\ell[k] \triangleq Z_\ell[k]/(BT)$  (nats/s/Hz). Then, the service process  $C_\ell[k]$  of the  $\ell$ th video layer is given by

$$C_\ell[k] = BT(t_\ell[k]R_\ell[k] + z_\ell[k]) \quad (\text{nats/frame}). \quad (4)$$

Clearly, the pre-drop strategy suppresses the growing speed of the queue for a more robust queueing behavior, but this strategy also causes data loss to all multicast receivers. As a result,  $z_\ell[k]$  needs to be determined by not only the CSI, but also the loss constraints (see Section III-D).

#### D. Loss Rate Constraint

Although a certain loss is usually tolerable for delay-sensitive services, the loss level cannot be arbitrarily high. Consequently, we require the loss rate of the  $\ell$ th video layer for each receiver to be limited lower than or equal to an application-dependent threshold, denoted by  $q_{\text{th}}^{(\ell)}$ . The loss rate of the  $\ell$ th video layer for the  $n$ th receiver, denoted by  $q_{\ell,n}$ , is defined as the ratio of the amount of data correctly received by this receiver to that of the data transmitted at this video layer. Data loss for unicast will be caused only by the pre-drop strategy, while data loss for multicast will be introduced by both pre-drop operation and heterogeneous channel fading across multicast receivers.

Since various efficient forward-error control (FEC) codes [16]-[18] at upper protocol layers were proposed and widely applied to multicast communications in wired/wireless networks, in our framework we suppose that FEC mechanisms are already employed at the upper protocol layers. The error-control redundancies added in the FEC codes at different video layers are inherently related among video layers and are jointly determined by the targeted video-delivery qualities at different video layers. Correspondingly, the tolerable loss-rate levels  $q_{\text{th}}^{(\ell)}$ 's for different video layers (indexed by  $\ell$ ) are jointly specified based on the video delivery quality requirements and the error control redundancy degrees across different video layers. Under this framework, we then mainly focus on how to use the minimum wireless resources with QoS guarantees to unicast/multicast multi-layer video over wireless channels.

### IV. STATISTICAL DELAY QOS GUARANTEES THROUGH EFFECTIVE BANDWIDTH/CAPACITY

#### A. Preliminary for Statistical Delay QoS Guarantees

The theory on statistical delay QoS guarantees [9]-[11] provides a powerful approach in analyzing the queueing behavior for time-varying arrival and/or service processes. Specifically, consider a stable queueing system with the stationary and

ergodic arrival and service processes. Asymptotic analyses [9] show that with sufficient conditions, the queue length process  $Q[k]$  converges to a random variable  $Q[\infty]$  in distribution such that

$$-\lim_{Q_{\text{th}} \rightarrow \infty} \frac{\log(\Pr\{Q[\infty] > Q_{\text{th}}\})}{Q_{\text{th}}} = \theta \quad (5)$$

for a certain  $\theta > 0$ , where  $Q_{\text{th}}$  is the queue-length bound. Moreover, the queue-length bound violation probability can be approximated by

$$\Pr\{Q > Q_{\text{th}}\} \approx e^{-\theta Q_{\text{th}}}, \quad (6)$$

where we remove the index  $[k]$  for  $Q[k]$  to simplify notations. In the above two equations,  $\theta$  is called *QoS exponent* and can be used to characterize delay QoS. The larger  $\theta$  corresponds to the more stringent QoS requirement, while the smaller  $\theta$  imposes the looser delay constraint. When the QoS metric of interest becomes delay-bound, the similar expressions can be obtained.

#### B. Preliminary for Effective Bandwidth/Capacity

Effective bandwidth [10] and effective capacity [11] are a pair of dual concepts. Given an arrival process  $A[k]$ , effective bandwidth of  $A[k]$ , denoted by  $\mathcal{A}(\theta)$  (nats/frame), is defined as the *minimum constant service rate* required to guarantee a specified QoS exponent  $\theta$ , i.e., Eqs. (5)-(6) are satisfied for the given  $\theta$ . In contrast, given a service process  $C[k]$ , effective capacity of  $C[k]$ , denoted by  $\mathcal{C}(\theta)$  (nats/frame), is defined as the *maximum constant arrival rate* which can be supported by  $C[k]$  subject to the specified QoS exponent  $\theta$ . Moreover, effective bandwidth [10] and effective capacity [11] can be expressed as  $\mathcal{A}(\theta) = \lim_{k \rightarrow \infty} (1/(k\theta)) \log(\mathbb{E}\{e^{\theta S_A[k]}\})$  and  $\mathcal{C}(\theta) = -\lim_{k \rightarrow \infty} (1/(\theta k)) \log(\mathbb{E}\{e^{-\theta S_C[k]}\})$ , respectively, where  $\mathbb{E}\{\cdot\}$  denotes the expectation,  $S_A[k] \triangleq \sum_{i=1}^k A[i]$ , and  $S_C[k] \triangleq \sum_{i=1}^k C[i]$ . In addition, for effective capacity and effective bandwidth scenarios, the delay-bound violation probability can be approximated [10], [11] as:

$$\Pr\{D > D_{\text{th}}\} \approx e^{-\theta \mathcal{C}(\theta) D_{\text{th}}}, \quad \text{for effective capacity;} \quad (7)$$

and

$$\Pr\{D > D_{\text{th}}\} \approx e^{-\theta \mathcal{A}(\theta) D_{\text{th}}}, \quad \text{for effective bandwidth.} \quad (8)$$

Equations (6), (7), and (8) are good approximations for relatively large  $Q_{\text{th}}$  and  $D_{\text{th}}$  as shown in [11], [22]. When  $Q_{\text{th}}$  and  $D_{\text{th}}$  are relatively small, the more accurate approximations expressions than Eqs. (6), (7), and (8) are given in [11], [22] as follows:

$$\begin{cases} \Pr\{Q > Q_{\text{th}}\} & \approx \varrho e^{-\theta Q_{\text{th}}}; \\ \Pr\{D > D_{\text{th}}\} & \approx \varrho e^{-\theta \mathcal{C}(\theta) D_{\text{th}}}; \\ \Pr\{D > D_{\text{th}}\} & \approx \varrho e^{-\theta \mathcal{A}(\theta) D_{\text{th}}}, \end{cases}$$

where  $\varrho$  denotes the probability that the queue is nonempty. These approximations are upper-bounded by the corresponding approximations given in Eqs. (6)-(8). Thus, directly using Eqs. (6)-(8) for the system design often guarantees more stringent QoS than the specified requirements. For wireless video transmissions, the delay bound  $D_{\text{th}}$  is typically hundreds of milliseconds (ms), which are thus much larger than the

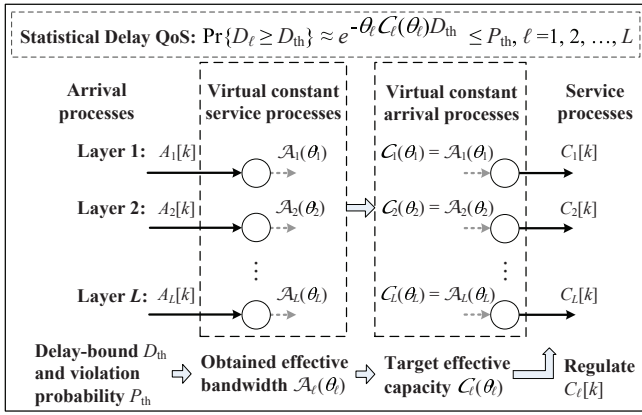


Fig. 2. Illustration of design procedures to guarantee statistical delay QoS by using effective capacity and effective bandwidth theories.

adaptive-transmission period scale (e.g., the physical-layer time-frame length) of the wireless transmission system, where the adaptive-transmission period typically varies from a few milliseconds (ms) to tens of milliseconds (ms). Therefore, Eqs. (6)-(8) are good approximating expressions in designing efficient wireless video-transmission schemes with statistical QoS guarantees.

### C. Design Procedures for Transmitting Layered Video with Statistical QoS Guarantees

For a dynamic queueing system, in order to guarantee the QoS exponent  $\theta$  given in Eqs. (5)-(6), the following equation need to be satisfied [10], [12]:

$$C(\theta) = \mathcal{A}(\theta). \quad (9)$$

Inspired by this property, the statistical delay QoS guarantees can be characterized through the arrival process and service process separately. As shown in Fig. 2, the queueing system for the  $\ell$ th video layer can be decomposed to two virtual queueing systems. The one on the left-hand side of Fig. 2 is composed by the true arrival process  $A_\ell[k]$  and one virtual constant-rate service process, the rate of which is equal to the effective bandwidth  $\mathcal{A}_\ell(\theta_\ell)$  of  $A_\ell[k]$ ; the right one consists of the true service process  $C_\ell[k]$  and one constant-rate virtual arrival process, the rate of which is equal to the effective capacity  $\bar{C}_\ell(\theta_\ell)$  of  $C_\ell[k]$ . Using the above concept, we develop the design procedures to provide statistical delay QoS guarantees for transmitting multi-layer video stream as shown in Fig. 3.

Among the procedures in Fig. 3, Steps 1 and 2 first identify the effective bandwidth  $\mathcal{A}_\ell(\theta_\ell)$  and QoS exponent  $\theta_\ell$  required to satisfy the delay-bound  $D_{th}$  and its violation probability  $P_{th}$ . Then, to satisfy the delay QoS in Eq. (1), we need to either satisfy Eq. (9) or guarantee that the effective capacity is larger than the effective bandwidth, which results in Steps 3 and 4.

The analytical expressions of effective bandwidth for many typical arrival processes, such as constant-rate process, autoregressive (AR) process, and Markovian process, can be found in [10]. Note that if  $A_\ell[k]$  is time varying, in Step 2  $\theta_\ell$  can be determined through Eq. (8). However, if  $A_\ell[k]$  is a

Step 1:	Determine effective bandwidth functions $\mathcal{A}_\ell(\theta)$ for the arrival processes $A_\ell[k]$ , $\ell = 1, 2, \dots, L$ .
Step 2:	Apply Eq. (7) or Eq. (8) to find the solution $\theta_\ell$ to the equation $\Pr\{D_\ell > D_{th}\} = P_{th}$ and get the corresponding effective bandwidth $\mathcal{A}_\ell(\theta_\ell)$ .
Step 3:	Set the target effective capacity $\bar{C}_\ell = \mathcal{A}_\ell(\theta_\ell)$ for the service processes of each video layer.
Step 4:	Jointly design adaptive service process $C_\ell[k]$ for each video layer, such that $C_\ell(\theta_\ell) \geq \bar{C}_\ell$ is satisfied while minimizing the total consumed wireless resources.

Fig. 3. The design procedures to provide statistical delay QoS guarantees for transmitting multi-layer video stream.

constant-rate process equal to  $\bar{A}_\ell$ , we have  $\mathcal{A}_\ell(\theta) = \bar{A}_\ell$  for all  $\theta$ , implying that the delay-bound violation probability of the virtual queueing system on the left-hand side of Fig. 2 is always equal to 0. Therefore, we cannot derive  $\theta_\ell$  directly through Eq. (8). In contrast, the QoS exponent  $\theta_\ell$  to guide the adaptive transmission needs to be determined by using Eq. (7) under the condition of  $C_\ell(\theta_\ell) = \mathcal{A}_\ell(\theta_\ell) = \bar{A}_\ell$ .

### V. UNICASTING MULTI-LAYER VIDEO STREAM

Assuming that the target effective capacity  $\{\bar{C}_\ell\}_{\ell=1}^L$  and QoS exponent  $\theta_\ell$  have been determined, we next focus on developing the optimal adaptive time-slot allocation and pre-drop strategy to satisfy the QoS requirements while minimizing the wireless-resource consumption. Unless otherwise mentioned, we drop the time-frame index  $[k]$  for the corresponding variables in the rest of this paper to simplify notations. Based on our previous work [13], if a stationary and ergodic service rate process  $C_\ell$  is uncorrelated across different time frames, we can write its effective capacity as follows:

$$C_\ell(\theta_\ell) = -\frac{1}{\theta_\ell} \log(\mathbb{E}\{e^{-\theta_\ell C_\ell}\}) \quad \text{nats/frame}. \quad (10)$$

Since the block-fading channel model described in Section II satisfies the time-uncorrelated condition, we can apply Eq. (10) for our framework to derive the adaptive unicast/multicast schemes with the statistical QoS guarantees.

#### A. Unicasting Layered Video Stream Without Loss Tolerance

We first consider the cases without loss tolerance for multi-layer video transmissions, i.e.,  $q_{th}^{(\ell)} = 0$  for all  $\ell$  and the pre-drop strategy will not be applied. Thus, we only need to focus on regulating the time-slot proportion  $\{t_\ell\}_{\ell=1}^L$  for each video layer. Following the design target and QoS constraints characterized in Sections III and IV, we derive the adaptive transmission strategy by solving the following optimization problem.

*Problem 1: P1*

$$\min_{\mathbf{t}} \left\{ \sum_{\ell=1}^L \mathbb{E}_\gamma \{t_\ell\} \right\} \quad (11)$$

$$\text{s.t.} \quad C_\ell(\theta_\ell) \geq \bar{C}_\ell, \quad \forall \ell, \quad (12)$$

$$\sum_{\ell=1}^L t_\ell - 1 \leq 0, \quad 0 \leq t_\ell \leq 1, \quad \forall \gamma, \quad (13)$$

where  $\mathbf{t} \triangleq (t_1, t_2, \dots, t_L)$  and  $\mathbb{E}_\gamma\{\cdot\}$  denotes the expectation over the random variable  $\gamma$ .

Using Eq. (10), the constraint in (12) can be equivalently rewritten as

$$\mathbb{E}_\gamma \left\{ e^{-\beta_\ell t_\ell \log(1+\gamma)} \right\} - V_\ell \leq 0, \quad (14)$$

where  $\beta_\ell \triangleq \theta_\ell TB$  is termed normalized QoS exponent and  $V_\ell \triangleq e^{-\theta_\ell \bar{C}_\ell}$ . It is not difficult to see: 1) the objective function in **P1** is convex over  $t$ ; 2) the functions on the left-hand side of all inequality constraints (Eqs. (13) and (14)) are convex over  $t$ . Therefore, **P1** is a convex problem [20] and the optimal solution can be obtained by using the Lagrangian method and the Karush-Kuhn-Tucker (KKT) conditions [20], which is summarized in Theorem 1.

*Theorem 1:* The optimal solution  $t^*$  to problem **P1**, if existing, is determined by

$$t_\ell^* = t_\ell(\gamma, \lambda_\ell^*, \psi_\gamma^*) \triangleq \left[ -\frac{\log\left(\frac{1+\psi_\gamma^*}{\beta_\ell \lambda_\ell^* \log(1+\gamma)}\right)}{\beta_\ell \log(1+\gamma)} \right]^+, \quad (15)$$

where  $[\cdot]^+ \triangleq \max\{\cdot, 0\}$ . The parameters  $\psi_\gamma^*$  and  $\{\lambda_\ell^*\}_{\ell=1}^L$  are the optimal Lagrangian multipliers associated with Eqs. (13) and (14), respectively. Given SNR  $\gamma$  in a fading state and  $\{\lambda_\ell^*\}_{\ell=1}^L$ , if

$$\sum_{\ell=1}^L t_\ell(\gamma, \lambda_\ell^*, 0) \geq 1 \quad (16)$$

holds,  $\psi_\gamma^*$  is the unique solution to

$$\sum_{\ell=1}^L t_\ell(\gamma, \lambda_\ell^*, \psi_\gamma^*) = 1, \quad \psi_\gamma^* \geq 0; \quad (17)$$

otherwise, we get

$$\psi_\gamma^* = 0. \quad (18)$$

Under the above strategy to determine  $t^*$  and  $\psi_\gamma^*$ , the optimal  $\{\lambda_\ell^*\}_{\ell=1}^L$  are selected to satisfy

$$\mathbb{E}_\gamma \left\{ e^{-\beta_\ell t_\ell^* \log(1+\gamma)} \right\} - V_\ell = 0, \quad \forall \ell. \quad (19)$$

*Proof:* We construct the Lagrangian function for **P1**, denoted by  $J$ , as follows:

$$J = \mathbb{E}_\gamma \left\{ \sum_{\ell=1}^L t_\ell \right\} + \mathbb{E}_\gamma \left\{ \psi_\gamma \left( \sum_{\ell=1}^L t_\ell - 1 \right) \right\} + \sum_{\ell=1}^L \lambda_\ell \left( \mathbb{E}_\gamma \left\{ e^{-\beta_\ell t_\ell \log(1+\gamma)} \right\} - V_\ell \right), \quad (20)$$

where  $\psi_\gamma \geq 0$  and  $\lambda_\ell \geq 0$ ,  $\ell = 1, 2, \dots, L$ , are Lagrangian multipliers associated with Eqs. (13) and (14), respectively. Then, the optimal  $t^*$  and Lagrangian multipliers of optimization problem **P1** satisfy the following KKT conditions [20]:

$$\begin{cases} \frac{\partial J}{\partial t_\ell} \Big|_{t_\ell=t_\ell^*} = 0, \forall \ell, \forall \gamma \\ \psi_\gamma^* \geq 0 \text{ and } \lambda_\ell^* \geq 0; \\ \psi_\gamma^* \left( \sum_{\ell=1}^L t_\ell^* - 1 \right) = 0, \forall \ell, \gamma; \\ \lambda_\ell^* \left( \mathbb{E}_\gamma \left\{ e^{-\beta_\ell t_\ell^* \log(1+\gamma)} \right\} - V_\ell \right) = 0, \forall \ell. \end{cases} \quad (21)$$

Taking the derivative of  $J$  with respect to (w.r.t.)  $t_\ell$ , we get

$$\frac{\partial J}{\partial t_\ell} = \left( 1 + \psi_\gamma - \lambda_\ell \beta_\ell (1 + \gamma)^{-\beta_\ell t_\ell} \log(1 + \gamma) \right) f_\Gamma(\gamma) d\gamma \quad (22)$$

where  $f_\Gamma(\gamma)$  is the probability density function (pdf) of  $\gamma$ . Plugging Eq. (22) into the first line of Eq. (21) and solving for  $t_\ell^*$  under the boundary condition  $t_\ell \geq 0$ , we get Eq. (15).

According to Eq. (15),  $t_\ell(\gamma, \lambda_\ell^*, \psi_\gamma)$  is a strictly decreasing function of  $\psi_\gamma$  for  $t_\ell(\gamma, \lambda_\ell^*, \psi_\gamma) > 0$ , and  $\psi_\gamma^* \rightarrow \infty$  leads to  $t_\ell^* \rightarrow 0$ . Therefore, if Eq. (16) holds, we can always find the unique  $\psi_\gamma^*$  to satisfy Eq. (18); otherwise, the inequality  $\sum_{\ell=1}^L t_\ell(\gamma, \lambda_\ell^*, \psi_\gamma) - 1 < 0$  follows for any  $\psi_\gamma \geq 0$ , implying  $\psi_\gamma^* = 0$  by applying the third line of Eq. (21). Through Eq. (15),  $\lambda_\ell^* = 0$  results in  $t_\ell = 0$  for all  $\gamma$ , and thus the constraint in Eq. (12) will be violated, implying an infeasible solution. Therefore,  $\lambda_\ell^*$  has to be positive. Then, to satisfy the fourth line of Eq. (21), Eq. (19) must hold and thus Theorem 1 follows. ■

Note that given  $\{\lambda_\ell^*\}_{\ell=1}^L$ ,  $\psi_\gamma^*$  is easy to solve because  $t_\ell(\gamma, \lambda_\ell^*, \psi_\gamma)$  is a decreasing function of  $\psi_\gamma$ . However, how to find  $\{\lambda_\ell^*\}_{\ell=1}^L$  is still unknown. Moreover, Theorem 1 does not state whether the optimal solution exists. Next, we discuss how to get  $\{\lambda_\ell^*\}_{\ell=1}^L$  and examine the existence of the optimal solution, which can be performed either off-line or on-line. Based on the optimization theory [19],[20], the Lagrangian dual problem to **P1** is given by

$$\max_{(\lambda, \psi_\gamma)} \left\{ \tilde{J}(\lambda, \psi_\gamma) \right\}, \quad (23)$$

where  $\lambda \triangleq (\lambda_1, \lambda_2, \dots, \lambda_L)$  and  $\tilde{J}(\lambda, \psi_\gamma)$  is the Lagrangian dual function defined by  $\tilde{J}(\lambda, \psi_\gamma) \triangleq \min_t \{J\} = J|_{t_\ell=t_\ell(\gamma, \lambda_\ell, \psi_\gamma)}$ . We can further convert Eq. (23) into  $\max_{(\lambda, \psi_\gamma)} \left\{ \tilde{J}(\lambda, \psi_\gamma) \right\} = \max_\lambda \left\{ \tilde{J}(\lambda, \psi_\gamma(\lambda)) \right\}$ , where  $\psi_\gamma(\lambda)$  denotes the maximizer of  $\tilde{J}(\lambda, \psi_\gamma)$  given  $\lambda$ . Moreover, we can obtain  $\psi_\gamma(\lambda)$  by using the same procedures as those used in determining  $\psi_\gamma^*$ , which are given by Eqs. (16)-(18) in Theorem 1.

Since problem **P1** is convex, there is no duality gap between **P1** and its dual problem given by Eq. (23) if the optimal solution exists. Thus, the optimal Lagrangian multipliers  $\{\lambda_\ell^*\}_{\ell=1}^L$  and  $\psi_\gamma^*$  to problem **P1** also maximize the objective function  $\tilde{J}(\lambda, \psi_\gamma)$  in Eq. (23). Consequently, we can obtain  $\{\lambda_\ell^*\}_{\ell=1}^L$  through maximizing  $\tilde{J}(\lambda, \psi_\gamma)$ . Following convex optimization theory [20],  $\tilde{J}(\lambda, \psi_\gamma(\lambda))$  is a concave function over  $\lambda$ , and thus we can track the optimal  $\lambda^*$  by using the subgradient method [21]:

$$\lambda_\ell := \lambda_\ell + \epsilon \left( \mathbb{E}_\gamma \left\{ e^{-\beta_\ell t_\ell \log(1+\gamma)} \right\} - V_\ell \right) \Big|_{t_\ell=t_\ell(\gamma, \lambda_\ell, \psi_\gamma(\lambda))} \quad (24)$$

where  $\epsilon$  is a positive real number close to 0, and  $(\mathbb{E}_\gamma \left\{ e^{-\beta_\ell t_\ell \log(1+\gamma)} \right\} - V_\ell)$  in the above equation is a subgradient of  $\tilde{J}(\lambda, \psi_\gamma(\lambda))$  w.r.t.  $\lambda_\ell$  [19] (see the definition of subgradient in Section VI-B). If the optimal solution to **P1** exists, the above iteration will converge to the optimal  $\lambda^*$  with properly selected  $\epsilon$  because of the concavity of  $\tilde{J}(\lambda, \psi_\gamma(\lambda))$ . Correspondingly,  $(\mathbb{E}_\gamma \left\{ e^{-\beta_\ell t_\ell \log(1+\gamma)} \right\} - V_\ell)$  will converge to 0. If the optimal solution to **P1** does not exist, we cannot

support such a statistical QoS requirement even we use up all time slots. Then,  $(\mathbb{E}_\gamma \{e^{-\beta_\ell t_\ell \log(1+\gamma)}\} - V_\ell)$  is always larger than 0 for some  $\ell$ . As a result,  $\lambda_\ell$  will approach infinity. So, if any  $\lambda_\ell$  does not converge and keeps increasing, we can conclude that the optimal solution does not exist and current wireless resources are not enough to support the specified statistical delay QoS for the incoming multi-layer video stream.

To find the optimal  $\lambda^*$ , we need the pdf of  $\gamma$ . In realistic systems, although the pdf of  $\gamma$  is usually unknown, we can still apply Eq. (24) to implement online tracking. In particular, the iterative update of Eq. (24) will be performed in each time frame. However, the expectation of  $e^{-\beta_\ell t_\ell \log(1+\gamma)}$  in Eq. (24) needs to be substituted by its estimation obtained based on the statistics from previous time frames. Denoting the estimation of  $\mathbb{E}_\gamma \{e^{-\beta_\ell t_\ell \log(1+\gamma)}\}$  in the  $k$ th time frame by  $\mathcal{S}_\ell[k]$ , we obtain  $\mathcal{S}_\ell[k+1]$  through a first-order autoregressive filter (low-pass filter) as follows:

$$\mathcal{S}_\ell[k+1] := (1-\alpha)\mathcal{S}_\ell[k] + \alpha e^{-\beta_\ell t_\ell [k+1] \log(1+\gamma[k+1])}, \quad (25)$$

where  $\ell = 1, 2, \dots, L$  and  $\alpha \in (0, 1)$  is a small positive number close to 0. If the optimal solution exists, the online tracking method converges with properly selected  $\alpha$  and  $\epsilon$ . Section VII will present some examples of tracking the optimal Lagrangian multipliers through simulations.

### B. Unicasting Layered Video Stream With Loss Tolerance

When  $q_{\text{th}}^{(\ell)} > 0$ , the transmission strategy becomes more complicated, but will use less wireless resources. After integrating the pre-drop strategy, the loss rate  $q_\ell$  of the  $\ell$ th layer is derived as

$$q_\ell = 1 - \frac{BT\mathbb{E}_\gamma \{t_\ell R_\ell\}}{\mathbb{E}_\gamma \{C_\ell\}} = 1 - \frac{\mathbb{E}_\gamma \{t_\ell R_\ell\}}{\mathbb{E}_\gamma \{t_\ell R_\ell + z_\ell\}}. \quad (26)$$

Next, we identify the adaptive transmission policy by solving optimization problem **P2**:

**Problem 2: P2**

$$\min_{(t, z)} \left\{ \sum_{\ell=1}^L \mathbb{E}_\gamma \{t_\ell\} \right\} \quad (27)$$

$$\text{s.t.: } \mathbb{E}_\gamma \left\{ e^{-\beta_\ell (z_\ell + t_\ell \log(1+\gamma))} \right\} - V_\ell \leq 0, \quad z_\ell \geq 0, \quad \forall \ell, \quad (28)$$

$$q_\ell \leq q_{\text{th}}^{(\ell)}, \quad \forall \ell \quad (29)$$

$$\sum_{\ell=1}^L t_\ell \leq 1, \quad 0 \leq t_\ell \leq 1, \quad \forall \gamma, \quad (30)$$

where  $\mathbf{z} \triangleq (z_1, z_2, \dots, z_L)$ .

Applying Eq. (26), we can rewrite Eq. (29) as follows:

$$\left(1 - q_{\text{th}}^{(\ell)}\right) \mathbb{E}_\gamma \{z_\ell\} - q_{\text{th}}^{(\ell)} \mathbb{E}_\gamma \{t_\ell \log(1+\gamma)\} \leq 0, \quad \forall \ell. \quad (31)$$

It is also not hard to prove that problem **P2** is still a convex problem and the Lagrangian method is still effective in finding the optimal solutions, which is summarized in Theorem 2.

**Theorem 2:** The optimal solution  $(t^*, z^*)$ , if existing, is expressed by a set of functions of  $\gamma$ ,  $\{\lambda_\ell\}_{\ell=1}^L$ ,  $\{\phi_\ell\}_{\ell=1}^L$ , and  $\psi_\gamma$  as follows::

$$t_\ell^* = t_\ell(\gamma, \lambda_\ell^*, \phi_\ell^*, \psi_\gamma^*), \quad z_\ell^* = z_\ell(\gamma, \lambda_\ell^*, \phi_\ell^*, \psi_\gamma^*), \quad \forall \ell, \quad (32)$$

where

$$t_\ell(\gamma, \lambda_\ell, \phi_\ell, \psi_\gamma) \triangleq \begin{cases} \infty, & \text{if } -\infty < \frac{1+\psi_\gamma}{\log(1+\gamma)} \leq \phi_\ell q_{\text{th}}^{(\ell)}; \\ \left[ -\frac{1}{\beta_\ell \log(1+\gamma)} \log \left( \frac{1+\psi_\gamma - q_{\text{th}}^{(\ell)} \phi_\ell \log(1+\gamma)}{\beta_\ell \lambda_\ell \log(1+\gamma)} \right) \right]^+, & \text{if } \phi_\ell q_{\text{th}}^{(\ell)} < \frac{1+\psi_\gamma}{\log(1+\gamma)} < \phi_\ell; \\ 0, & \text{if } \phi_\ell \leq \frac{1+\psi_\gamma}{\log(1+\gamma)} < \infty, \end{cases} \quad (33)$$

and

$$z_\ell(\gamma, \lambda_\ell, \phi_\ell, \psi_\gamma) \triangleq \begin{cases} 0, & \text{if } -\infty < \frac{1+\psi_\gamma}{\log(1+\gamma)} < \phi_\ell; \\ \left[ -\frac{1}{\beta_\ell} \log \left( \frac{\phi_\ell (1 - q_{\text{th}}^{(\ell)})}{\beta_\ell \lambda_\ell} \right) \right]^+, & \text{if } \phi_\ell \leq \frac{1+\psi_\gamma}{\log(1+\gamma)} < \infty. \end{cases} \quad (34)$$

Given  $\gamma$ ,  $\{\lambda_\ell^*\}_{\ell=1}^L$ , and  $\{\phi_\ell^*\}_{\ell=1}^L$ , if  $\sum_{\ell=1}^L t_\ell(\gamma, \lambda_\ell^*, \phi_\ell^*, 0) \geq 1$ ,  $\psi_\gamma^*$  is selected such that the equation  $\sum_{\ell=1}^L t_\ell(\gamma, \lambda_\ell^*, \phi_\ell^*, \psi_\gamma^*) = 1$  holds; otherwise, we have  $\psi_\gamma^* = 0$ . The optimal  $\{\lambda_\ell^*\}_{\ell=1}^L$  and  $\{\phi_\ell^*\}_{\ell=1}^L$  need to be jointly selected such that “=” holds in both Eqs. (28) and (31).

*Proof:* The proof of Theorem 2 can be readily obtained by using the standard Lagrangian-multiplier based method and KKT conditions, and thus is omitted due to lack of space. The detailed proof of Theorem 2 is provided online in [24]. ■

In order to search for the optimal Lagrangian multipliers and check the existence of the optimal solution, we can also design the adaptive tracking method similar to problem **P1**.

## VI. QOS GUARANTEES FOR MULTICASTING LAYERED-VIDEO STREAM

We consider the multicast scenario in this section. If no loss is tolerated, the transmission rate in each time frame is limited by the worst-case SNR among all multicast receivers. Thus, the system throughput will be degraded very quickly as the multicast group size increases. Therefore, we mainly focus on the multicast scenario with loss tolerance.

### A. Problem Formulation for Multicast Scenario

Under the multicast rate-adaptation strategy given in Section III, the loss rate of the  $n$ th receiver at the  $\ell$ th video layer becomes

$$q_{\ell, n} = 1 - \frac{\mathbb{E}_\gamma \left\{ t_\ell \sum_{i=1}^N t_{\ell, i} \log(1+\gamma_i) \delta_{\gamma_n \geq \gamma_i} \right\}}{\mathbb{E}_\gamma \{z_\ell + t_\ell R_\ell\}}, \quad (35)$$

where  $\mathbb{E}_\gamma \{\cdot\}$  denotes the expectation over all fading states of the random vector variable  $\gamma$ ,  $\delta_{\gamma_n \geq \gamma_i}$  is the indication function (for a given statement  $\varpi$ ,  $\delta(\varpi) = 1$  if  $\varpi$  is true, and  $\delta(\varpi) = 0$  otherwise), and  $R_\ell = \sum_{i=1}^N t_{\ell, i} \log(1+\gamma_i)$  is the total normalized transmission rate in a time frame (see Eq. (3)). Accordingly, the following loss-rate constraint needs to be satisfied for each multicast receiver at every video layer, which is specified by the inequality as follows:

$$q_{\ell, n} \leq q_{\text{th}}^{(\ell)}, \quad \forall n, \forall \ell. \quad (36)$$



To simplify the derivations, we first use a *relaxed* constraint to replace Eq. (36) by:

$$q_{\ell,0} \triangleq \frac{1}{N} \sum_{n=1}^N q_{\ell,n} \leq q_{\text{th}}^{(\ell)}, \quad \forall \ell \quad (37)$$

where  $q_{\ell,0}$  is called *group loss rate* (average loss rate over receivers) at the  $\ell$ th video layer. We will show later that the optimal adaptation policy derived under the group-loss-rate constraint given in Eq. (37) does not violate Eq. (36), and thus is also optimal under the original loss-rate constraint given by Eq. (36). Plugging Eq. (35) into Eq. (37), we have

$$q_{\ell,0} = 1 - \frac{\mathbb{E}_{\gamma} \left\{ t_{\ell} \sum_{i=1}^N t_{\ell,i} m_i \log(1 + \gamma_i) \right\}}{N \mathbb{E}_{\gamma} \{ z_{\ell} + t_{\ell} R_{\ell} \}}, \quad (38)$$

where  $m_i$  is the number of receivers with SNR higher than or equal to  $\gamma_i$ . In addition, it is clear that  $R_{\ell}$  falls in the following range:

$$R_{\ell} \in [R_{\pi(N)}, R_{\pi(1)}], \quad (39)$$

where  $R_{\pi(N)} \triangleq \min_{1 \leq n \leq N} \{\log(1 + \gamma_n)\}$  and  $R_{\pi(1)} \triangleq \max_{1 \leq n \leq N} \{\log(1 + \gamma_n)\}$ . Note that when we attempt to use a normalized transmission rate equal to  $R_{\ell}$  in a time frame, there are many different choices for  $\{t_{\ell,n}\}_{n=1}^N$  to get the same  $R_{\ell}$ . In order to minimize the loss for the entire multicast group, among all these choices we need to select the one which maximizes the numerator of the second term on the right-hand side of Eq. (38), which represents the sum rate of data correctly received by each multicast receiver. Accordingly, we define

$$\begin{aligned} \tilde{g}_s(R_{\ell}) &\triangleq \max_{\mathbf{t}_{\ell}: \sum_{i=1}^N t_{\ell,i} = 1} \left\{ \sum_{i=1}^N t_{\ell,i} m_i \log(1 + \gamma_i) \right\} \\ \text{s.t.} &: \sum_{i=1}^N t_{\ell,i} \log(1 + \gamma_i) = R_{\ell} \end{aligned} \quad (40)$$

where  $\mathbf{t}_{\ell} \triangleq (t_{\ell,1}, t_{\ell,2}, \dots, t_{\ell,N})$ . Therefore,  $\tilde{g}_s(R_{\ell})$  denotes the maximum sum of achieved rates over all multicast receivers under the given normalized transmission rate  $R_{\ell}$ . In our previous work [15], where we studied the tradeoff between the average throughput and loss rate for single-layer multicast, we showed that  $\tilde{g}_s(R_{\ell})$  can be derived through the concept of convex hull [20]. Using the properties of convex hull, in [15] we proved that  $\tilde{g}_s(R_{\ell})$  is a *continuous, piecewise linear*, and *concave function* over  $R_{\ell}$ . Thus, we can obtain  $\tilde{g}_s(R_{\ell})$  as follows:

$$\tilde{g}_s(R_{\ell}) = \begin{cases} \tilde{g}_s(r_i) + \eta_i(R_{\ell} - r_i), & \text{if } R_{\ell} \in [r_i, r_{i-1}), \\ & 2 \leq i \leq N; \\ \tilde{g}_s(r_2) + \eta_2(r_1 - r_2), & \text{if } R_{\ell} = r_1, \end{cases} \quad (41)$$

where  $R_{\pi(1)} = r_1 > r_2 > \dots > r_N = R_{\pi(N)}$ . Fig. 4 depicts an example for the function  $\tilde{g}_s(R_{\ell})$ . As shown in Fig. 4, within each interval  $[r_i, r_{i-1})$ ,  $\tilde{g}_s(R_{\ell})$  is a linear function of  $R_{\ell}$  with the slope equal to  $\eta_i$ , and  $(N - 1)$  is equal to the number of such intervals. Note that  $\{(r_i, \tilde{g}_s(r_i))\}_{i=1}^N$  are actually the vertices on the upper boundary of the convex hull of the 2-dimensional point set  $\{(\log(1 + \gamma_i), m_i \log(1 + \gamma_i))\}_{i=1}^N$  (see [15]). For the complete procedures to identify  $\tilde{g}_s(R_{\ell})$

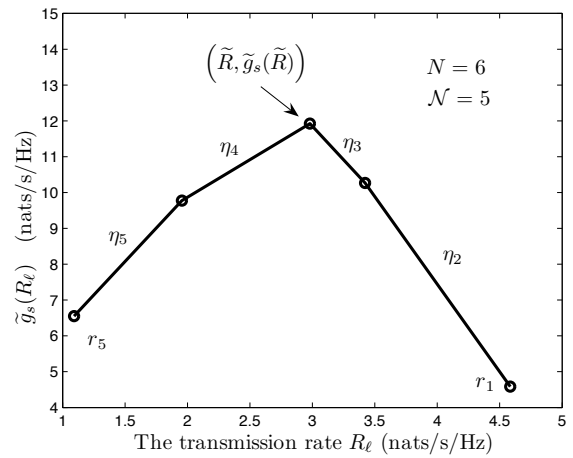


Fig. 4. An example for the function  $\tilde{g}_s(R_{\ell})$  against  $R_{\ell}$ , where  $N = 6$  and  $\gamma = (1.98, 6.07, 18.71, 29.63, 45.43, 96.93)$ .

and the corresponding time slot allocation policy, please refer to our previous work [15]. The above discussions imply that we need to consider only the transmission policies yielding  $\tilde{g}_s(R_{\ell})$ , because only these policies can minimize the total loss for the entire multicast group. Moreover, Eqs. (40)-(41) suggest that we can focus on regulating the scalar  $R_{\ell}$  instead of the  $N$ -dimension time-proportion vector  $\mathbf{t}_{\ell}$ . After  $R_{\ell}$  is determined, we can use Eq. (40) to obtain  $\mathbf{t}_{\ell}$ .

Following previous analyses in this section, we formulate problem **P3** to derive the adaptation policy for multi-layer video multicast as follows:

**Problem 3: P3**

$$\min_{(\mathbf{t}, \mathbf{R}, \mathbf{z})} \left\{ \sum_{\ell=1}^L \mathbb{E}_{\gamma} \{ t_{\ell} \} \right\} \quad (42)$$

$$\text{s.t.} \quad \mathbb{E}_{\gamma} \left\{ e^{-\beta_{\ell}(z_{\ell} + t_{\ell} R_{\ell})} \right\} - V_{\ell} \leq 0, \quad z_{\ell} \geq 0, \quad \forall \ell, \quad (43)$$

$$N \left( 1 - q_{\text{th}}^{(\ell)} \right) \mathbb{E}_{\gamma} \{ t_{\ell} R_{\ell} + z_{\ell} \} - \mathbb{E}_{\gamma} \{ t_{\ell} \tilde{g}_s(R_{\ell}) \} \leq 0, \quad \forall \ell, \quad (44)$$

$$\sum_{\ell=1}^L t_{\ell} - 1 \leq 0, \quad 0 \leq t_{\ell} \leq 1, \quad \forall \ell, \quad (45)$$

where  $\mathbf{R} \triangleq (R_1, R_2, \dots, R_L)$  and Eq. (44) is the group-loss-rate constraint (equivalent to Eq. (37)) for the policies corresponding to  $\tilde{g}_s(R_{\ell})$ .

## B. Derivation of the Optimal Solution for Multicast Video

Notice that Problem **P3** is not convex, because the functions on the left-hand side of Eqs. (43) and (44) are not convex over  $(\mathbf{t}, \mathbf{R}, \mathbf{z})$ . However, we show in the following that the optimal solution can still be obtained through Lagrangian dual problem.

### B.1 Lagrangian Characterization of Problem P3

The Lagrangian function of **P3**, denoted by  $W$ , is constructed as

$$W = \mathbb{E}_{\gamma} \{ w \} \quad (46)$$



where

$$w \triangleq \sum_{\ell=1}^L t_\ell + \psi_\gamma \left( \sum_{\ell=1}^L t_\ell - 1 \right) + \sum_{\ell=1}^L \lambda_\ell \left( e^{-\beta_\ell(z_\ell + t_\ell R_\ell)} - V_\ell \right) + \sum_{\ell=1}^L \phi_\ell \left( N \left( 1 - q_{\text{th}}^{(\ell)} \right) (t_\ell R_\ell + z_\ell) - t_\ell \tilde{g}_s(R_\ell) \right). \quad (47)$$

In Eq. (47),  $\psi_\gamma \geq 0$ ,  $\phi_\ell \geq 0$ , and  $\lambda_\ell \geq 0$  are the Lagrangian multipliers associated with the constraints given by Eqs. (45), (44), and (43), respectively. The Lagrangian dual function, denoted by  $U$ , is then determined by

$$U(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma) \triangleq \min_{(\boldsymbol{t}, \boldsymbol{z}, \boldsymbol{R})} \{W\} = \mathbb{E}_\gamma \{u(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma)\}, \quad (48)$$

where  $\boldsymbol{\phi} \triangleq (\phi_1, \phi_2, \dots, \phi_L)$  and

$$u(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma) \triangleq \min_{(\boldsymbol{t}, \boldsymbol{z}, \boldsymbol{R})} \{w\}. \quad (49)$$

It is clear that  $u(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma)$  is a concave function over  $(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma)$ , and so is  $U(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma)$ . Moreover, the Lagrange dual problem is defined as:

$$\mathbf{P3-Dual} : U^* \triangleq U(\boldsymbol{\lambda}^*, \boldsymbol{\phi}^*, \psi_\gamma^*) = \max_{(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma)} \{U(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma)\}, \quad (50)$$

where  $(\boldsymbol{\lambda}^*, \boldsymbol{\phi}^*, \psi_\gamma^*)$  is the maximizer. We then solve for the optimal adaptation strategy to **P3** through the dual problem. If the optimal solution to **P3** exists, we will show later that there is no duality gap between the primal problem **P3** and the dual problem **P3-Dual**. As a result, the optimal solution to **P1** must minimize the Lagrangian function  $W$  under the optimal Lagrangian multipliers  $(\boldsymbol{\lambda}^*, \boldsymbol{\phi}^*, \psi_\gamma^*)$ .

## B.2 Derivation of the Lagrangian Dual Function

$$U(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma) = \mathbb{E}_\gamma \{u(\boldsymbol{\lambda}, \boldsymbol{\phi}, \psi_\gamma)\}$$

Since  $\tilde{g}_s(R_\ell)$  is nondifferentiable at some  $R_\ell$  (as shown in Fig. 4),  $w$  is also nondifferentiable at some  $R_\ell$ . Alternatively, we need to use the subgradient and subdifferential [19] instead of gradient to derive the minimizer to Eq. (49), which is denoted by  $(\boldsymbol{t}^*, \boldsymbol{z}^*, \boldsymbol{R}^*)$ .

*Definition 1:* Consider a convex function  $h : \mathcal{D} \rightarrow \mathbb{R}$ , where  $\mathbb{R}$  denotes the set of real numbers and  $\mathcal{D} \subset \mathbb{R}^n$  is a convex set. Then, an  $n \times 1$  vector  $\boldsymbol{\xi}$ , for  $\boldsymbol{\xi} \in \mathbb{R}^n$ , is called a *subgradient* [19] at  $\mathbf{d} \in \mathcal{D}$  if  $h(\mathbf{d}') \geq h(\mathbf{d}) + \boldsymbol{\xi}^T(\mathbf{d}' - \mathbf{d})$  for all  $\mathbf{d}' \in \mathcal{D}$ , where  $(\cdot)^T$  represents the transpose. The collection of subgradients at  $\mathbf{d}$  form a set called the *subdifferential* [19] of  $h(\cdot)$  at  $\mathbf{d}$ , denoted by  $\partial h(\mathbf{d})$ . If  $h(\cdot)$  is differentiable at  $\mathbf{d}$ , the subgradient at  $\mathbf{d}$  is unique and becomes the gradient. Moreover, the sufficient and necessary condition that  $\mathbf{d}^*$  minimizes  $h(\mathbf{d})$  is that  $\mathbf{0} \in \partial h(\mathbf{d}^*)$ . When  $h(\cdot)$  is a *concave* function, the subgradient and subdifferential at  $h(\mathbf{d})$  is defined in the similar way as the convex case, except that the required inequality becomes  $h(\mathbf{d}') \leq h(\mathbf{d}) + \boldsymbol{\xi}^T(\mathbf{d}' - \mathbf{d})$ .

Then, using Eq. (40) and Definition 1, we obtain

$$\partial \tilde{g}_s(r) = \begin{cases} [\eta_N, \infty], & \text{if } r = r_N; \\ \{\eta_i\}, & \text{if } r_i < r < r_{i-1}, 2 \leq i \leq N; \\ [\eta_i, \eta_{i+1}], & \text{if } r = r_i, 2 \leq i \leq N-1; \\ [-\infty, \eta_1], & \text{if } r = r_1. \end{cases} \quad (51)$$

Applying Eq. (47) into Definition 1, we get the subdifferential of  $w$  w.r.t.  $R_\ell$ , denoted by  $\partial w_{R_\ell}$ , as:

$$\partial w_{R_\ell} = \left\{ y \mid y = t_\ell \left( \phi_\ell N \left( 1 - q_{\text{th}}^{(\ell)} \right) - \lambda_\ell \beta_\ell e^{-\beta_\ell(z_\ell + t_\ell R_\ell)} - \phi_\ell x \right), \forall x \in \partial \tilde{g}_s(R_\ell) \right\}, \forall \ell, \gamma. \quad (52)$$

It is clear that  $w$  is differentiable w.r.t.  $t_\ell$  and  $R_\ell$ . Taking the derivative of  $w$  w.r.t.  $t_\ell$  and  $z_\ell$ , respectively, we get

$$\frac{\partial w}{\partial t_\ell} = 1 + \psi_\gamma - \lambda_\ell \beta_\ell R_\ell e^{-\beta_\ell(z_\ell + t_\ell R_\ell)} - \phi_\ell \tilde{g}_s(R_\ell) + \phi_\ell N \left( 1 - q_{\text{th}}^{(\ell)} \right) R_\ell, \forall \ell, \gamma; \quad (53)$$

$$\frac{\partial w}{\partial z_\ell} = \phi_\ell N \left( 1 - q_{\text{th}}^{(\ell)} \right) - \lambda_\ell \beta_\ell e^{-\beta_\ell(z_\ell + t_\ell R_\ell)}, \forall \ell, \gamma. \quad (54)$$

Clearly, the minimization of  $w$  can be performed separately for each video layer. Now consider the  $\ell$ th video layer. Since the function  $w$  is not convex over the 3-tuple  $(t_\ell, z_\ell, R_\ell)$ , the equations  $\partial w / \partial t_\ell = 0$ ,  $\partial w / \partial z_\ell = 0$ , and  $0 \in \partial w_{R_\ell}$  are only the necessary conditions for  $(t_\ell^*, z_\ell^*, R_\ell^*)$ . However, if  $t_\ell^*$  is given,  $w$  becomes a convex function over the 2-tuple  $(z_\ell, R_\ell)$ . Using this property, we can decompose the minimization of  $w$  into several easier sub-problems. Applying the above principle, we discuss the cases with the fixed  $t_\ell^* = 0$  and  $t_\ell^* > 0$ , respectively, as follows.

1)  $t_\ell^* = 0$ : The variable  $R_\ell$  vanishes in Eq. (47) when  $t_\ell = 0$ . Then, we only need to find the minimizer  $z_\ell^*$ . By solving  $\partial w / \partial z_\ell = 0$  under the condition  $z_\ell \geq 0$ , we get

$$z_\ell^* = \left[ -\frac{1}{\beta_\ell} \log \left( \frac{\phi_\ell N (1 - q_{\text{th}}^{(\ell)})}{\lambda_\ell \beta_\ell} \right) \right]^+. \quad (55)$$

2)  $t_\ell^* > 0$ : We jointly solve  $\partial w / \partial z_\ell = 0$  and  $0 \in \partial w_{R_\ell}$  under the condition  $z_\ell \geq 0$  and  $R_\ell \geq 0$ , and then get the minimizer, which is summarized in Eqs. (56)-(57) as follows:

if  $\widehat{R}_\ell > \widetilde{R}$ , then

$$\begin{cases} R_\ell^* = \widetilde{R}; \\ z_\ell^* = \left[ -\frac{1}{\beta_\ell} \log \left( \frac{\phi_\ell N (1 - q_{\text{th}}^{(\ell)})}{\lambda_\ell \beta_\ell} \right) - t_\ell^* R_\ell^* \right]^+; \end{cases} \quad (56)$$

if  $\widehat{R}_\ell \leq \widetilde{R}$ , then

$$\begin{cases} R_\ell^* = \widehat{R}_\ell; \\ z_\ell^* = 0, \end{cases} \quad (57)$$

where  $\widehat{R}_\ell$  is the unique solution to

$$0 \in (\partial w_{R_\ell})|_{z_\ell=0} \quad (58)$$

under the given  $t_\ell$ , and

$$\widetilde{R} \triangleq \arg \max_r \{ \tilde{g}_s(r) \}. \quad (59)$$

The detailed derivations for Eqs. (56) and (57) are omitted due to lack of space, but are provided online in [24]. Note that  $\widetilde{R}$  depends only on  $\tilde{g}_s(r)$ , but not on  $t_\ell^*$ . Then, through Eqs. (56)-(57), we can see that with  $t_\ell^* > 0$ , the minimizer must satisfy either  $z_\ell^* = 0$  or  $R_\ell^* = \widetilde{R}$ . Further note that the above results provide not only the mathematical convenience, but also the insightful observations for the adaptive multicast transmission. For multicast, zero loss can be achieved only

when setting transmission rate  $R_\ell = R_{\pi(N)}$ , which is determined by the worst-case SNR over all multicast receivers. The higher the transmission rate or the higher the drop rate we use, the higher the loss rate we get (observed from Eq. (38)). Therefore, when not violating the statistical delay QoS guarantees, we need to choose the transmission rate  $R_\ell$  and the drop rate  $z_\ell$  as small as possible. Moreover,  $R_\ell$  and  $z_\ell$  need to be jointly derived for performance optimization. Following the above strategies, we first derive  $\hat{R}_\ell$  which optimizes the system performance (equivalently, minimizes the Lagrangian function) given the zero drop rate. However, if  $\hat{R}_\ell > \bar{R}$ , we can see that the achieved sum rate  $\tilde{g}_s(R_\ell)$  over all multicast receivers under  $R_\ell = \hat{R}_\ell$  decreases when  $\hat{R}_\ell$  increases, as depicted in Fig. 4. When this happens, we need to set  $R_\ell = \bar{R}$  and apply the nonzero drop rate to avoid the degradation of  $\tilde{g}_s(R_\ell)$  while supporting the satisfied service rate.

Based on the above results for  $t_\ell^* = 0$  and  $t_\ell^* > 0$ , the minimizer  $(t_\ell^*, R_\ell^*, z_\ell^*)$  at the  $\ell$ th video layer must fall into one of the following three Sub-domains:

$$\begin{cases} \text{Sub-domain 1: } t_\ell = 0, R_\ell \geq 0, z_\ell \geq 0; \\ \text{Sub-domain 2: } t_\ell \geq 0, R_\ell = \bar{R}, z_\ell \geq 0; \\ \text{Sub-domain 3: } t_\ell \geq 0, R_\ell \geq 0, z_\ell = 0. \end{cases} \quad (60)$$

In Eq. (60), Sub-domain 1 is associated with the case of  $t_\ell^* = 0$ . For the case with  $t_\ell^* > 0$ , Sub-domains 2 and 3 correspond to the conditions  $\hat{R}_\ell \geq \bar{R}$  and  $\hat{R}_\ell < \bar{R}$ , respectively. In order to get the minimizer  $(t_\ell^*, R_\ell^*, z_\ell^*)$  of  $w$ , we can first find the minimizer within each Sub-domain, which is denoted by  $(t_\ell^{(j)}, R_\ell^{(j)}, z_\ell^{(j)})$ ,  $j = 1, 2, 3$ . After identifying the minimizers of each Sub-domain,  $(t_\ell^*, R_\ell^*, z_\ell^*)$  can then be obtained through

$$(t_\ell^*, R_\ell^*, z_\ell^*) = (t_\ell^{(j^*)}, R_\ell^{(j^*)}, z_\ell^{(j^*)}) \quad \forall \ell, \quad (61)$$

where

$$j^* = \arg \min_{j=1,2,3} \left\{ w \Big|_{(t_\ell, R_\ell, z_\ell) = (t_\ell^{(j)}, R_\ell^{(j)}, z_\ell^{(j)})} \right\}.$$

In Sub-domains 1, 2, and 3, the variables  $t_\ell$ ,  $R_\ell$ , and  $z_\ell$  are fixed, respectively, implying that there are only two optimization variables in each Sub-domain. Therefore, the minimization problem within each Sub-domain becomes tractable. For Sub-domain 1, the minimizer  $(0, R_\ell^{(1)}, z_\ell^{(1)})$  is given in Eq. (55). For Sub-domain 2, since  $R_\ell$  is fixed, deriving the minimizer  $(t_\ell^{(2)}, R_\ell^{(2)}, z_\ell^{(2)})$  is equivalent to solving a convex problem. For Sub-domain 3, the optimization problem can be readily solved by applying the piecewise linear property of  $\tilde{g}(R_\ell)$ . The detailed derivations for  $(t_\ell^{(2)}, R_\ell^{(2)}, z_\ell^{(2)})$  and  $(t_\ell^{(3)}, R_\ell^{(3)}, z_\ell^{(3)})$  are omitted due to lack of space, but are provided online in [24].

### B.3. The Optimal Solution to **P3**.

In Section VI-B.2, we have obtained the minimizer  $(t^*, z^*, \mathbf{R}^*)$  for  $w$ . Then, based on the optimization theory [19], the necessary and sufficient conditions for zero duality gap are as follows: there exists the feasible policy

$(t^*, z^*, \mathbf{R}^*)|_{\{\phi_\gamma = \phi_\gamma^*, \lambda_\ell = \lambda_\ell^*, \phi_\ell = \phi_\ell^*, \forall \ell, \gamma\}}$  such that

$$\begin{cases} \psi_\gamma^* \left( \sum_{\ell=1}^L t_\ell^* - 1 \right) = 0, \quad \forall \ell, \gamma; \\ \lambda_\ell^* \left( \mathbb{E}_\gamma \left\{ e^{-\beta_\ell(z_\ell^* + t_\ell^* R_\ell^*)} \right\} - V_\ell \right) = 0, \quad \forall \ell; \\ \phi_\ell^* \mathbb{E}_\gamma \left\{ N \left( 1 - q_{\text{th}}^{(\ell)} \right) \left( t_\ell^* R_\ell^* + z_\ell^* \right) - t_\ell^* \tilde{g}_s(R_\ell^*) \right\} = 0, \quad \forall \ell; \\ \psi_\gamma^* \geq 0, \lambda_\ell^* \geq 0, \phi_\ell^* \geq 0. \end{cases} \quad (62)$$

Then, the optimal policy to **P3** is given by

$$(t^*, z^*, \mathbf{R}^*)|_{\{\phi_\gamma = \phi_\gamma^*, \lambda_\ell = \lambda_\ell^*, \phi_\ell = \phi_\ell^*, \forall \ell, \gamma\}}. \quad (63)$$

We can solve for the optimal Lagrangian multipliers by using the similar arguments in the proof of Theorem 1. Specifically, for each channel realization if  $\sum_{\ell=1}^L t_\ell^* |_{\psi_\gamma=0} \leq 1$ , we have  $\psi_\gamma^* = 0$ ; otherwise,  $\psi_\gamma^*$  is the unique solution to  $\sum_{\ell=1}^L t_\ell^* = 1$ . Moreover,  $\{\phi_\ell^*\}_{\ell=1}^L$  and  $\{\lambda_\ell^*\}_{\ell=1}^L$  will be selected such that “=” holds for constraints given in Eqs. (43)-(44). Furthermore, we can design the adaptive tracking method similar to problem **P1** to examine the existence of the optimal solution and find the optimal Lagrangian multipliers.

Note that under the above optimal solution, different  $\{\gamma_n\}_{n=1}^N$ , which have the same ordered permutation, will generate the same function  $\tilde{g}_s(R_\ell)$  defined by Eq. (40) and thus the same adaptation policy. Then, since  $\gamma_n$ 's are i.i.d. (as assumed in Section II), this policy will benefit all receivers evenly, implying  $q_{\ell,0} = q_{\ell,1} = q_{\ell,2} = \dots = q_{\ell,N} = q_{\text{th}}^{(\ell)}$ . Therefore, the original loss-rate constraint is not violated for all multicast receivers. Moreover, since the group-loss-rate constraint given by Eq. (44) in problem **P3** is a relaxed version of the original loss-rate constraint for each multicast receiver (given by Eq. (36)), the optimal solution to problem **P3** is also optimal even if we replace Eq. (44) by using the original loss-rate constraint given in Eq. (36).

## VII. SIMULATION EVALUATIONS

We use simulation experiments to evaluate the performances of our proposed optimal adaptive transmission schemes and to investigate the impact of QoS requirements on resource allocations. Note that the metric “delay” investigated/simulated in simulations represents the queueing delay, as addressed previously in Section III-B for the framework of this paper. In simulations, we set the signal bandwidth  $B$  and time-frame length  $T$  equal to  $2 \times 10^5$  Hz and 10 ms, respectively. The arrival video stream includes two layers, both of which have constant arrival rates, where  $A_1[k] = 250$  Kbps and  $A_2[k] = 150$  Kbps. Then, the effective bandwidths of  $A_1[k]$  and  $A_2[k]$  are determined by  $\mathcal{A}_1(\theta_1) = 1.733 \times 10^3$  nats/frame and  $\mathcal{A}_2(\theta_2) = 1.040 \times 10^3$  nats/frame, respectively. The values of  $\theta_1$  and  $\theta_2$  can be derived from solving Eq. (7), depending on the QoS requirements specified by  $D_{\text{th}}$  and  $P_{\text{th}}$ ,  $\ell = 1, 2, \dots, L$ . The wireless channel follows the Rayleigh fading model and we denote the average SNR by  $\bar{\gamma}$ . Fig. 5 plots the iterative on-line tracking of the optimal Lagrangian multipliers  $\lambda_\ell^*$ 's based on the method used in Section V-A with  $\epsilon = 0.01$  and  $\alpha = 0.02$ . As shown in Fig. 5, the Lagrangian multipliers quickly converge to the optimal value and oscillate slightly within the small dynamic ranges, which demonstrates the effectiveness of our tracking method.

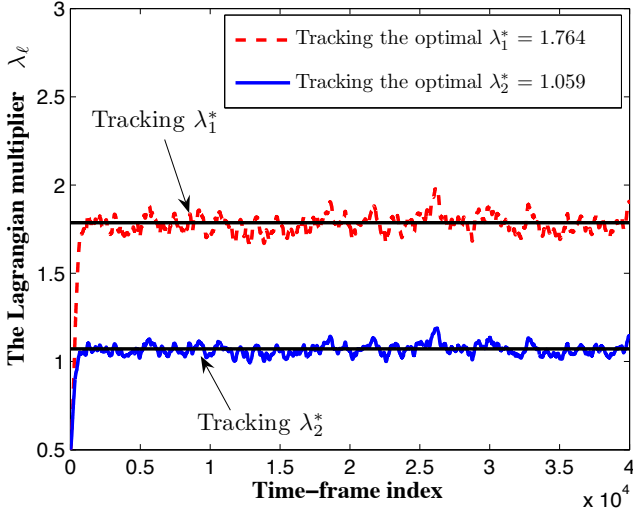


Fig. 5. Illustration of tracking the optimal Lagrangian multiplier  $\lambda_\ell^*$  for unicast with zero loss, where the average SNR is  $\bar{\gamma} = 10$  dB, the required delay bound is  $P_{th} = 10^{-4}$ , and the required threshold for the delay-bound violation probability is  $D_{th} = 250$  ms.

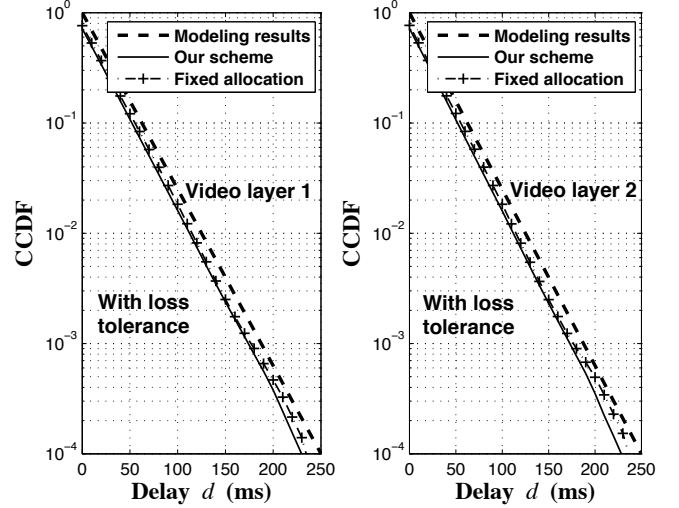


Fig. 7. The complementary cumulative distribution function  $\Pr\{D_\ell > d\}$  of the queuing delay for the unicast scenario with loss tolerance for video layer 1 and video layer 2, respectively, where  $\bar{\gamma} = 15$  dB,  $P_{th} = 10^{-4}$ ,  $D_{th} = 250$  ms, and  $(q_{th}^{(1)}, q_{th}^{(2)}) = (0.01, 0.02)$ .

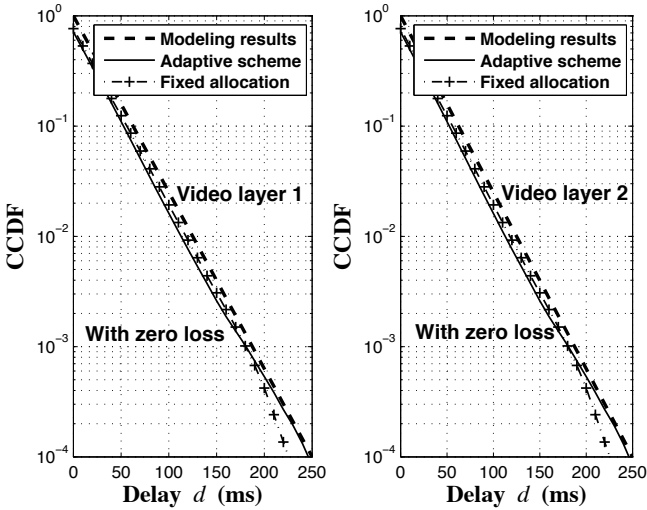


Fig. 6. The complementary cumulative distribution function (CCDF)  $\Pr\{D_\ell > d\}$  of the queuing delay for the unicast scenario with zero loss for video layer 1 and video layer 2, respectively, where  $\bar{\gamma} = 15$  dB,  $P_{th} = 10^{-4}$ , and the required delay-bound is  $D_{th} = 250$  ms.

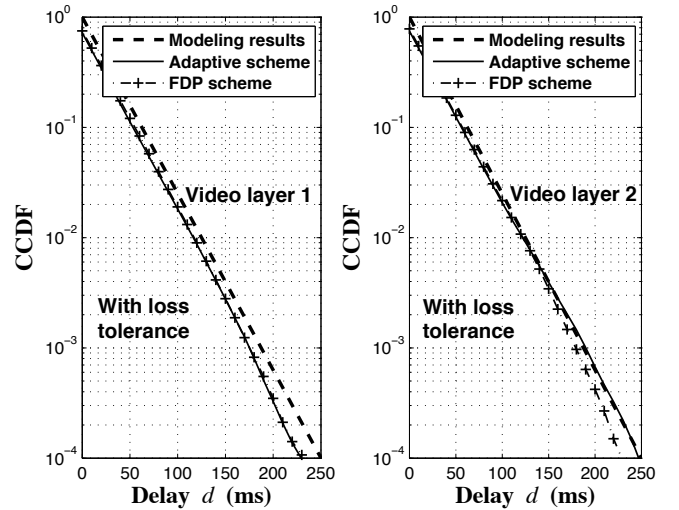


Fig. 8. The complementary cumulative distribution function  $\Pr\{D_\ell > d\}$  of the queuing delay for multicast scenario with loss tolerance for video layer 1 and video layer 2, respectively, where  $N = 20$  receivers,  $\bar{\gamma} = 20$  dB,  $(q_{th}^{(1)}, q_{th}^{(2)}) = (0.01, 0.05)$ ,  $P_{th} = 10^{-4}$ , and  $D_{th} = 250$  ms.

We also investigate some straightforward time-slot allocation schemes as the baseline schemes for comparative analyses. We will compare the average resource consumption between our derived optimal schemes and these baseline schemes under the same QoS satisfactions.

1) *Fixed time-slot allocation for unicast without loss*: This scheme uses constant time-slot length  $t_\ell[k] = \bar{t}_\ell$ ,  $k = 1, 2, \dots$ , in all fading states. The normalized transmission rate is set to  $R_\ell = \log(1 + \gamma)$  nats/s/Hz for the  $\ell$ th video layer. The parameters  $\bar{t}_\ell$ ,  $\ell = 1, 2$ , are selected such that the effective capacity  $\bar{C}_\ell(\theta_\ell)$  of the  $\ell$ th video layer's service process is just equal to  $\bar{C}_\ell$ :

$$C_\ell(\theta_\ell) = -\frac{1}{\theta_\ell} \log \left( \mathbb{E} \left\{ e^{-\theta_\ell \bar{t}_\ell B T \log(1+\gamma)} \right\} \right) = \bar{C}_\ell, \quad (64)$$

where  $\bar{C}_\ell = \mathcal{A}_\ell(\theta_\ell)$  and  $\ell = 1, 2, \dots, L$ . We can obtain the unique  $\bar{t}_\ell$ 's by numerically solving the above equation. If we get  $\sum_{\ell=1}^L \bar{t}_\ell > 1$ , this scheme cannot guarantee the QoS requirements under current channel conditions, even using up all time-slot resources.

2) *Fixed time-slot allocation for unicast with loss tolerance*: This scheme uses both the constant time-slot length  $t_\ell[k] = \bar{t}_\ell$  and the constant per-drop rate  $z_\ell[k] = \bar{z}_\ell$  in all fading states. The normalized transmission rate is also set to  $R_\ell = \log(1 + \gamma)$  nats/s/Hz for the  $\ell$ th video layer. The parameters  $\bar{t}_\ell$  and  $\bar{z}_\ell$  can be obtained by solving

$$C_\ell(\theta_\ell) = -\frac{1}{\theta_\ell} \log \left( \mathbb{E} \left\{ e^{-\theta_\ell (\bar{t}_\ell B T \log(1+\gamma) + B T \bar{z}_\ell)} \right\} \right) = \bar{C}_\ell \quad (65)$$

and  $q_\ell = q_{th}^{(\ell)}$  for all video layers, where  $\ell = 1, 2, \dots, L$ .

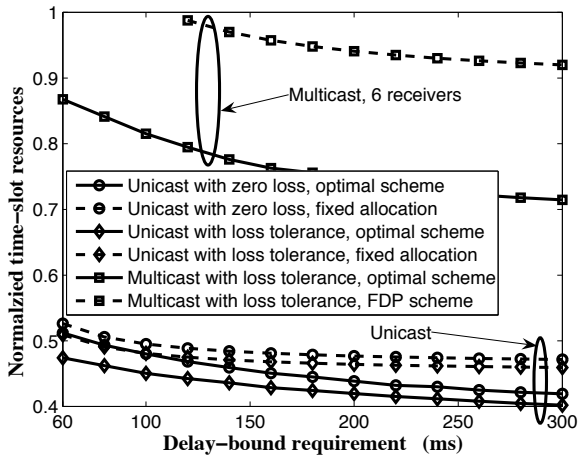


Fig. 9. The normalized time-slot resource consumption  $\mathbb{E}_\gamma \left\{ \sum_{\ell=1}^L t_\ell \right\}$  versus delay-bound requirement  $D_{\text{th}}$ , where  $\bar{\gamma} = 15$  dB,  $P_{\text{th}} = 10^{-4}$ , and  $(q_{\text{th}}^{(1)}, q_{\text{th}}^{(2)}) = (0.01, 0.05)$ .

3) *Fixed dominating position scheme for multicast with loss tolerance*: The fixed dominating position (FDP) scheme always sets  $R_\ell = \log(1 + \gamma_{\pi(i_\ell)})$  nats/s/Hz, where  $\gamma_{\pi(i)}$  denotes the  $i$ th largest instantaneous SNR among all multicast receivers. The index  $i_\ell$  is fixed at  $i_\ell = \left\lceil N \left( 1 - q_{\text{th}}^{(\ell)} \right) \right\rceil$  such that the loss-rate QoS is not violated. Moreover, the FDP scheme also adopts the constant time-slot length  $t_\ell[k] = \bar{t}_\ell$  and the constant per-drop rate  $z_\ell[k] = \bar{z}_\ell$ , which can be obtained by solving

$$C_\ell(\theta_\ell) = -\frac{1}{\theta_\ell} \log \left( \mathbb{E} \left\{ e^{-\theta_\ell (\bar{t}_\ell B T \log(1 + \gamma_{\pi(i_\ell)}) + B T \bar{z}_\ell)} \right\} \right) = \bar{C}_\ell \quad (66)$$

and  $q_{\ell,0} = q_{\text{th}}^{(\ell)}$  for all video layers, where  $\ell = 1, 2, \dots, L$ .

Figures 6, 7, and 8 depict the complementary cumulative distribution function (CCDF) of the queueing delay, i.e., the probability  $\Pr\{D_\ell > d\}$  given a threshold  $d$ , for unicast with zero loss, unicast with loss tolerance, and multicast with loss tolerance, respectively. We can observe from Figs. 6-8 that the CCDF's of all schemes agree well with the modeling results (see Eqs. (7)-(8)) at each video layer, where the delay-bound violation probability decreases exponentially against the delay bound. Moreover, for the required delay bound  $D_{\text{th}} = 250$  ms, the violation probability of all schemes can be upper-bounded by the targeted  $P_{\text{th}} = 10^{-4}$  for each video layer, which demonstrates the validity of all schemes in terms of statistical delay-QoS guarantees. Having shown that all schemes can meet the same QoS requirements, we then focus on the performance of the average time-slot consumption.

Figures 9 and 10 illustrate the impact of delay-bound  $D_{\text{th}}$  and its violation probability  $P_{\text{th}}$  on the resource consumption, respectively. As shown in Figs. 9 and 10, either smaller  $D_{\text{th}}$  or  $P_{\text{th}}$  will cause more resource consumption, which is expected because the smaller  $D_{\text{th}}$  or  $P_{\text{th}}$  implies the more stringent delay QoS requirement. Moreover, under various QoS conditions, our proposed optimal schemes always use much less wireless resources than the baseline schemes. For multicast services, our derived optimal scheme consumes at

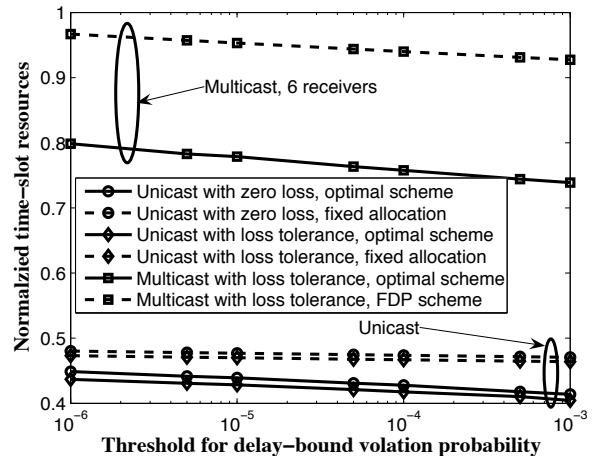


Fig. 10. The normalized time-slot resource consumption  $\mathbb{E}_\gamma \left\{ \sum_{\ell=1}^L t_\ell \right\}$  versus the threshold  $P_{\text{th}}$  for the delay-bound violation probability, where  $\bar{\gamma} = 15$  dB,  $D_{\text{th}} = 250$  ms, and  $(q_{\text{th}}^{(1)}, q_{\text{th}}^{(2)}) = (0.01, 0.02)$ .

least 15% of total resources less than the FDP scheme. For unicast services, more resources can be saved by using the optimal scheme when delay QoS becomes looser, while less resources are saved under the more stringent QoS constraints.

After demonstrating the superiority of our proposed optimal schemes over the baseline schemes, Fig. 11(a) plots the average time-slot consumption versus the average SNR for multi-layer video unicast and multicast. We observe from Fig. 11(a) that under the same channel conditions, video unicast uses much less time-slot resources than multicast. When the average SNR is relatively low around 7-8 dB, video unicast does not need to consume all available time-slot resources. In contrast, video multicast almost uses up *all* resources even with  $\bar{\gamma} = 13.5$  dB for the 6-receiver case. For the 10-receiver case, the average SNR needs to be larger than 15 dB to provide the QoS-guaranteed multicast services. The above observations reflect the key challenges on wireless multicast. In wireless broadcast channels, since all receivers can hear the sender, it is ideal that only one copy of data is transmitted such that sizable resources can be saved. However, due to the heterogeneous fading channels across multicast receivers, the transmission rate has to be limited within the relatively low range to avoid too much data loss for receivers with poorer instantaneous channel qualities. As a result, more time-slot resources are consumed to meet the QoS requirements. In addition, more multicast receivers result in more resource consumption, as depicted in Fig. 11(a). But note that although the wireless multicast faces many challenges, it still uses much less wireless resources than the strategy which uses multiple unicast links to implement wireless multicast. For example, if using multiple unicast links to implement multicast, we need the time-slot resources at least  $N$  times as much as the resource consumption for a unicast link. Clearly, for environments simulated in Fig. 11(a), even with  $\bar{\gamma} = 18$  dB, we still do not have enough resources for such a unicast-based multicast scheme with just 6 receivers. Fig. 11(b) shows the resource consumption for each video layer. We can see that video layer 1 requires more resources

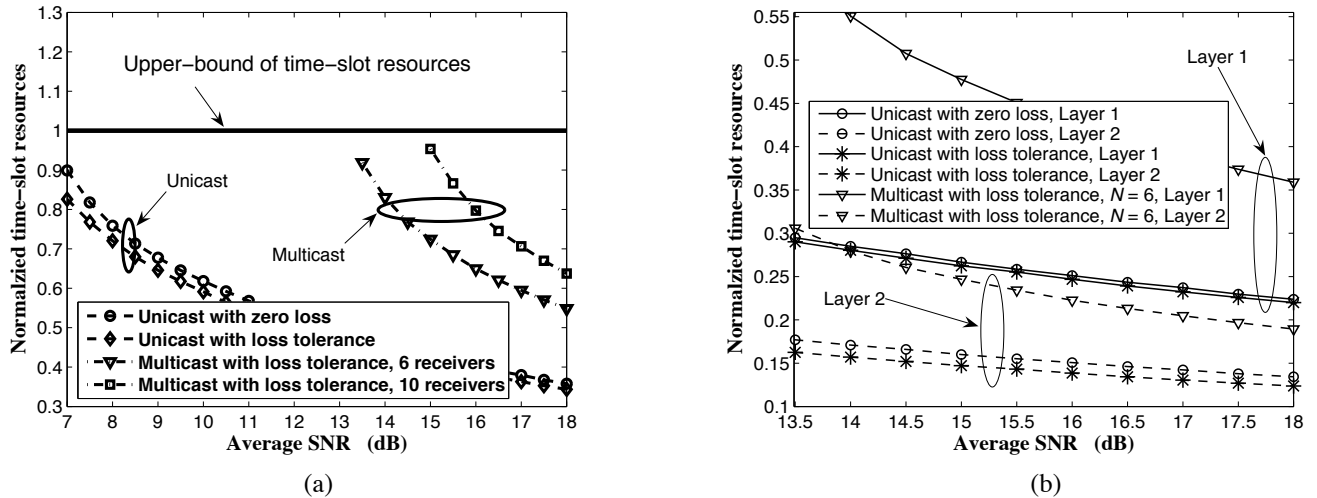


Fig. 11. (a) The normalized time-slot resource consumption  $\mathbb{E}_\gamma \left\{ \sum_{\ell=1}^L t_\ell \right\}$  versus  $\bar{\gamma}$ , where  $P_{\text{th}} = 10^{-4}$ ,  $D_{\text{th}} = 250$  ms,  $(q_{\text{th}}^{(1)}, q_{\text{th}}^{(2)}) = (0.01, 0.05)$ . (b)  $\mathbb{E}_\gamma \{t_\ell\}$  for each video layer versus  $\bar{\gamma}$ , where  $P_{\text{th}} = 10^{-4}$ ,  $D_{\text{th}} = 250$  ms,  $(q_{\text{th}}^{(1)}, q_{\text{th}}^{(2)}) = (0.01, 0.05)$ .

than video layer 2, which is because in our settings the traffic load of layer 1 is higher and the loss-rate QoS of layer 1 is more stringent. Furthermore, in multicast the difference of resource consumption between the two video layers is larger than the difference in unicast.

Figure 12 shows the impact from loss-rate constraints on video multicast. As shown in Fig. 12, even slightly increasing  $q_{\text{th}}^{(\ell)}$  can significantly reduce the total consumed wireless resources. This is because the higher loss-tolerance level will enable larger multicast transmission rate and thus consume less time-slot resources. The above observations suggest that there exists a tradeoff between loss-rate control at the physical layer and error recovery at the upper protocol layers. As mentioned previously, the loss-rate  $q_{\text{th}}^{(\ell)}$  depends largely on the capability of erasure-correction codes used at upper protocol layers, especially for multicast services. Thus, Fig. 12 suggests that using more redundancy for forward error-control at upper protocol layers can effectively decrease the total wireless-resource consumption with QoS guarantees, while enabling the repair of more data losses at the physical layer.

## VIII. CONCLUSIONS

We proposed a framework to model the wireless transmission of multi-layer video stream with statistical delay QoS guarantees. A separate queue is maintained for each video layer and the same statistical delay QoS-requirement needs to be satisfied by all video layers, where the statistical delay QoS is characterized by the delay-bound and its corresponding violation probability through the effective capacity bandwidth/capacity theory. Under the proposed framework, we derived a set of optimal rate adaptation and time-slot allocation schemes for video unicast/multicast with and/or without loss tolerance, which minimizes the time-slot resource consumption. We also conducted extensive simulation experiments to demonstrate the impact of statistical QoS provisionings on wireless resource allocations by using our derived optimal adaptive transmission schemes.

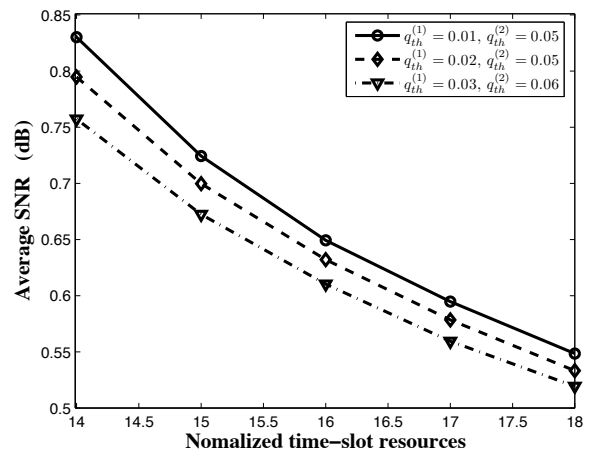


Fig. 12. The normalized time-slot resource consumption  $\mathbb{E}_\gamma \left\{ \sum_{\ell=1}^L t_\ell \right\}$  versus  $\bar{\gamma}$  for multicast with  $N = 6$ ,  $P_{\text{th}} = 10^{-4}$ , and  $D_{\text{th}} = 250$  ms.

## REFERENCES

- [1] J. Shin, J.-W. Kim, C.-C.J. Kuo, "Quality-of-service mapping mechanism for packet video indifferiated services network", *IEEE Trans. Multimedia*, vol. 3, no. 2, pp. 291-231, Jun. 2001.
- [2] H. Schwarz, D. Marpe, and T. Wiegand, "Overview of the scalable video coding extension of the H.264/AVC standard," *IEEE Trans. Circuits and Systems for Video Technology*, vol. 17, no. 9, pp. 1103-1120, Sep. 2007.
- [3] S. McCanne, M. Vetterli, and V. Jacobson, "Low-complexity video coding for receiver-driven layered multicast," *IEEE J. Sel. Areas Commun.*, vol. 15, no. 6, pp. 983-1001, Aug. 1997.
- [4] J. Villalon, P. Cuenca, L. L. Orozco-Barbosa, Y. Seok, and T. Turletti, "Cross-layer architecture for adaptive video multicast streaming over multirate wireless LANs," *IEEE J. Sel. Area. Commun.*, vol. 25, no. 4, pp. 699-711, May 2007.
- [5] X. Zhang and K. G. Shin, D. Saha, and D. Kandlur, "Scalable flow control for multicast ABR services in ATM networks," *IEEE/ACM Trans. on Netw.*, vol. 10, no. 1, pp. 67-85, Feb. 2002.
- [6] X. Zhang and K. G. Shin, "Delay analysis of feedback-synchronization signaling for multicast flow control," *IEEE/ACM Trans. on Netw.*, vol. 11, no.3, pp. 436-460, Jun. 2003.
- [7] X. Zhang and K. G. Shin, "Markov-chain modeling for multicast signaling delay analysis," *IEEE/ACM Trans. on Netw.*, vol. 12, no. 4, pp. 667-680, Aug. 2004.

- [8] X. Zhang and Q. Du, "Cross-layer modeling for QoS-driven multimedia multicast/broadcast over fading channels," *IEEE Commun. Mag.*, vol. 45, no. 8, pp. 62-70, Aug. 2007.
- [9] C.-S. Chang, "Stability, queue length, and delay of deterministic and stochastic queueing networks," *IEEE Trans. Auto. Control*, vol. 39, no. 5, pp. 913-931, May 1994.
- [10] C.-S. Chang, *Performance Guarantees in Communication Networks*, Springer-Verlag London, 2000.
- [11] D. Wu and R. Negi, "Effective capacity: A wireless link model for support of Quality of Service," *IEEE Trans. Wireless Commun.*, vol. 2, no. 4, pp. 630-643, Jul. 2003.
- [12] X. Zhang, J. Tang, H.-H. Chen, S. Ci, and M. Guizni, "Cross-layer-based modeling for quality of service guarantees in mobile wireless networks," *IEEE Commun. Mag.*, pp. 100-106, Jan. 2006.
- [13] J. Tang and X. Zhang, "Quality-of-service driven power and rate adaptation over wireless links," *IEEE Trans. Wireless Commun.*, vol. 6, no. 8, pp. 3058-3068, Aug. 2007.
- [14] Q. Du and X. Zhang, "Effective capacity of superposition coding based mobile multicast in wireless networks," in *Proc. IEEE International Conference on Commun., (ICC'09)*, Dresden, Germany, Jun. 2009.
- [15] Q. Du and X. Zhang, "Cross-layer design based rate control for mobile multicast in cellular networks," in *Proc. IEEE GLOBECOM 2007*, Washington DC, USA, Nov. 26-30, 2007, pp. 5180-5184.
- [16] J. Nonenmacher, E. Biersack, and D. Towsley, "Parity-based loss recovery for reliable multicast transmission," *IEEE/ACM Trans. Networking*, vol. 6, no. 4, pp. 349-361, Aug. 1998.
- [17] J. W. Byers, M. Luby, and M. Mitzenmacher, "A digital Fountain approach to asynchronous reliable multicast," *IEEE J. Select. Areas Commun.*, vol. 20, no. 8, pp. 1528-1540, Oct. 2002.
- [18] X. Zhang and Q. Du, "Adaptive Low-Complexity Erasure-Correcting Code Based Protocols for QoS-Driven Mobile Multicast Services Over Wireless Networks," *IEEE Trans. Veh. Technol.*, vol. 55, no. 5, pp. 1633-1647, Sep. 2006.
- [19] M. S. Bazaraa, H. D. Sherali, and C. M. Shetty, *Nonlinear Programming: Theory and Algorithms*, 3rd ed., John Wiley & Sons, Inc., 2006.
- [20] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [21] Stephen Boyd, Lin Xiao, and Almir Mutapcic, *Subgradient Methods*, [http://www.stanford.edu/class/ee392o/subgrad\\_method.pdf](http://www.stanford.edu/class/ee392o/subgrad_method.pdf).
- [22] G. L. Choudhury, D. M. Lucantoni, and W. Whitt, "Squeezing the most out of ATM," *IEEE Trans. Commun.*, vol. 44, pp. 203-217, Feb. 1996.
- [23] T.H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*, 2nd ed. Cambridge, Mass. and New York: MIT Press and McGraw-Hill, 2001.
- [24] Q. Du and X. Zhang, "Statistical QoS provisionings for wireless unicast/multicast of multi-layer video streams," *Networking and Information Systems Labs., Dept. Electr. and Comput. Eng., Texas A&M Univ., College Station, Tech. Rep.* [Online.] Available: [http://www.ece.tamu.edu/~xzhang/papers/multi\\_layer\\_video.pdf](http://www.ece.tamu.edu/~xzhang/papers/multi_layer_video.pdf).

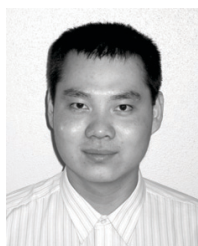


**Xi Zhang** [S'89-SM'98] received the B.S. and M.S. degrees from Xidian University, Xi'an, China, the M.S. degree from Lehigh University, Bethlehem, PA, all in electrical engineering and computer science, and the Ph.D. degree in electrical engineering and computer science (Electrical Engineering-Systems) from The University of Michigan, Ann Arbor.

He is currently an Associate Professor and the Founding Director of the Networking and Information Systems Laboratory, Department of Electrical and Computer Engineering, Texas A&M University, College Station. He was an Assistant Professor and the Founding Director of the Division of Computer Systems Engineering, Department of Electrical Engineering and Computer Science, Beijing Information Technology Engineering Institute, China, from 1984 to 1989. He was a Research Fellow with the School of Electrical Engineering, University of Technology, Sydney, Australia, and the Department of Electrical and Computer Engineering, James Cook University, Australia, under a Fellowship from the Chinese National Commission of Education. He worked as a Summer Intern with the Networks and Distributed Systems Research Department, AT&T Bell Laboratories, Murray Hills, NJ, and with AT&T Laboratories Research, Florham Park, NJ, in 1997. He has published more than 170 research papers in the areas of wireless networks and communications systems, mobile computing, network protocol design and modeling, statistical communications, random signal processing, information theory, and control theory and systems.

Prof. Zhang received the U.S. National Science Foundation CAREER Award in 2004 for his research in the areas of mobile wireless and multicast networking and systems. He received the Best Paper Awards in the IEEE Globecom 2009 and the IEEE Globecom 2007, respectively. He also received the TEES Select Young Faculty Award for Excellence in Research Performance from the Dwight Look College of Engineering at Texas A&M University, College Station, in 2006. He is currently serving as an Editor for the IEEE TRANSACTIONS ON COMMUNICATIONS, an Editor for the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, an Associate Editor for the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, a Guest Editor for the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS for the special issue on "wireless video transmissions", an Associate Editor for the IEEE COMMUNICATIONS LETTERS, a Guest Editor for the IEEE WIRELESS COMMUNICATIONS MAGAZINE for the special issue on "next generation of CDMA versus OFDMA for 4G wireless applications", an Editor for the JOHN WILEY'S JOURNAL ON WIRELESS COMMUNICATIONS AND MOBILE COMPUTING, an Editor for the JOURNAL OF COMPUTER SYSTEMS, NETWORKING, AND COMMUNICATIONS, an Associate Editor for the JOHN WILEY'S JOURNAL ON SECURITY AND COMMUNICATIONS NETWORKS, an Area Editor for the ELSEVIER JOURNAL ON COMPUTER COMMUNICATIONS, and a Guest Editor for JOHN WILEY'S JOURNAL ON WIRELESS COMMUNICATIONS AND MOBILE COMPUTING for the special issue on "next generation wireless communications and mobile computing". He has frequently served as the Panelist on the U.S. National Science Foundation Research-Proposal Panels. He is serving or has served as the Technical Program Committee (TPC) Chair for IEEE Globecom 2011, TPC Vice-Chair for IEEE INFOCOM 2010, TPC Co-Chair for IEEE INFOCOM 2009 Mini-Conference, TPC Co-Chair for IEEE Globecom 2008 - Wireless Communications Symposium, TPC Co-Chair for the IEEE ICC 2008 - Information and Network Security Symposium, Symposium Chair for IEEE/ACM International Cross-Layer Optimized Wireless Networks Symposium 2006, 2007, and 2008, respectively, the TPC Chair for IEEE/ACM IWCMC 2006, 2007, and 2008, respectively, the Demo/Poster Chair for IEEE INFOCOM 2008, the Student Travel Grants Co-Chair for IEEE INFOCOM 2007, the General Chair for ACM QShine 2010, the Panel Co-Chair for IEEE ICCCN 2007, the Poster Chair for IEEE/ACM MSWiM 2007 and IEEE QShine 2006, Executive Committee Co-Chair for QShine, the Publicity Chair for IEEE/ACM QShine 2007 and IEEE WirelessCom 2005, and the Panelist on the Cross-Layer Optimized Wireless Networks and Multimedia Communications at IEEE ICCCN 2007 and WiFi-Hotspots/WLAN and QoS Panel at IEEE QShine 2004. He has served as the TPC members for more than 70 IEEE/ACM conferences, including IEEE INFOCOM, IEEE Globecom, IEEE ICC, IEEE WCNC, IEEE VTC, IEEE/ACM QShine, IEEE WoWMoM, IEEE ICCCN, etc.

Prof. Zhang is a Senior Member of the IEEE and a Member of the Association for Computing Machinery (ACM).



**Qinghe Du** [S'09] received B.S. and M.S. from Xian Jiaotong University, and is currently working towards to the Ph.D. degree under supervising of Prof. Xi Zhang in Networking and Information Systems Laboratory, Department of Electrical and Computer Engineering at Texas A&M University.

His research interests include mobile wireless communications and networks with emphasis on mobile multicast, statistical QoS provisioning, QoS-driven resource allocations, cognitive radio techniques, and cross-layer design over wireless networks.

His work co-authored with his Ph.D. advisor Prof. Xi Zhang received the Best Paper Award in the IEEE Globecom 2007 for the paper "Cross-Layer Design Based Rate Control for Mobile Multicast in Cellular Networks". He has published multiple papers in IEEE Transactions, IEEE J-SAC, IEEE Comm Magazine, IEEE INFOCOM, IEEE Globecom, IEEE ICC, etc.