

Quality-Aware Streaming and Scheduling for Device-to-Device Video Delivery

Joongheon Kim, *Member, IEEE*, Giuseppe Caire, *Fellow, IEEE*, and Andreas F. Molisch, *Fellow, IEEE*

Abstract—On-demand video streaming is becoming a killer application for wireless networks. Recent information-theoretic results have shown that a combination of caching on the users' devices and device-to-device (D2D) communications yields throughput scalability for very dense networks, which represent critical bottlenecks for conventional cellular and wireless local area network (WLAN) technologies. In this paper, we consider the implementation of such caching D2D systems where each device pre-caches a subset of video files from a library, and users requesting a file that is not already in their library obtain it from neighboring devices through D2D communication. We develop centralized and distributed algorithms for the delivery phase, encompassing a link scheduling and a streaming component. The centralized scheduling is based on the max-weighted independent set (MWIS) principle and uses message-passing to determine max-weight independent sets. The distributed scheduling is based on a variant of the FlashLinQ link scheduling algorithm, enhanced by introducing video-streaming specific weights. In both cases, the streaming component is based on a quality-aware stochastic optimization approach, reminiscent of current Dynamic Adaptive Streaming over HTTP (DASH) technology, for which users sequentially request video “chunks” by choosing adaptively their quality level. The streaming and the scheduling components are coupled by the length of the users' request queues. Through extensive system simulation, the proposed approaches are shown to provide sizeable gains with respect to baseline schemes formed by the concatenation of off-the-shelf FlashLinQ with proportional fair link scheduling and DASH at the application layer.

Index Terms—Adaptive streaming, device-to-device, quality awareness, scheduling, video delivery.

I. INTRODUCTION

ACCORDING to recent predictions of the Cisco Visual Networking Index (VNI) [1], the sum of all forms of video will constitute 80%–90% of global consumer data traffic by 2017, and the traffic from wireless and mobile devices will exceed the traffic from wired devices by 2016. Therefore, efficient video-aware network algorithms for wireless networks are

of highest importance [2]–[23]. It has been shown recently [2], [17], [19]–[24] that the throughput for delivery of wireless video files can be greatly enhanced by device-to-device (D2D) communications, where direct links between pairs of user devices can be set up without requiring to go through a central base station.¹ In particular, these works propose systems where each device caches independently, according to a certain optimal distribution, a subset of popular video files. When a user needs a file not already present in its own cache, it obtains it from one of its neighbors through a spectrally efficient, short-range D2D link. As user density increases, the aggregate storage capacity of the D2D network increases linearly with the number of users, while the average communication distance decreases (and the spatial reuse increases). For these reasons, D2D networks for video delivery are scalable, such that the aggregate demand and throughput increase linearly with user density.

The most common way users consume video is by streaming, i.e., after a pre-buffering time typically much shorter than the duration of the video file, playback is started while, at the same time, the rest of the file is progressively downloaded. More specifically, the video file is divided into “chunks” such that while the already-received chunks are played in sequence, the later chunks are transmitted. A transmission algorithm for such a system consists of two components: 1) a scheduling algorithm that determines which D2D pairs are allowed to transmit at a given time, and 2) a streaming algorithm that determines adaptively from which device each chunk should be requested and at which quality level. These two components are obviously coupled.

Currently, the most well-known D2D scheduling protocol in both industry and academia is *FlashLinQ* [25], [26]. It is a distributed algorithm that schedules D2D links according to their priorities, such that the higher-priority links do not suffer from significant interference of possibly scheduled lower-priority links. Theoretically, it can guarantee the maximum number of activated D2D links as analyzed by the theory of stochastic geometry [27]. However, *FlashLinQ* does not incorporate naturally a video quality-aware mechanism, and therefore its suitability for D2D on-demand video streaming remains open.

As far as adaptive video streaming is concerned, a number of protocols have been suggested, considering various characteristics of video streaming.

For example, cloud-based video streaming protocols are proposed in [28] and [29], channel-aware streaming algorithms are discussed in [30], new architectural concepts are presented in [31], rate-distortion-theory-based (or quality-aware)

Manuscript received June 01, 2014; revised March 22, 2015; accepted June 29, 2015; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor S. Ramasubramanian. Date of publication July 31, 2015; date of current version August 16, 2016. This work was supported in part by a grant from Intel Labs, Cisco Systems, and Verizon Wireless in the framework of the Video Aware Wireless Networks (VAWN) Research Program and the National Science Foundation (NSF) under Grants CCF-1423140 and CNS-1457340.

J. Kim was with the Department of Computer Science, University of Southern California, Los Angeles, CA 90089 USA. He is now with Intel Corporation, Santa Clara, CA 95054 USA (e-mail: joongheok@usc.edu).

G. Caire is with the Department of Electrical Engineering, University of Southern California, Los Angeles, CA 90089 USA, and also with the Department of Electrical Engineering and Computer Science, Technical University of Berlin, Berlin 10623, Germany (e-mail: caire@usc.edu).

A. F. Molisch is with the Department of Electrical Engineering, University of Southern California, Los Angeles, CA 90089 USA (e-mail: molisch@usc.edu).

Digital Object Identifier 10.1109/TNET.2015.2452272

¹Strictly speaking, this definition refers to single-hop D2D networks, which are the topic of this paper. In the following, we will simply refer to them as “D2D” as there is no possibility of confusion.

streaming is addressed in [32] and [33], and resource-aware streaming algorithms are mentioned in [34].

The key postulate of this paper is that D2D scheduling and streaming are *coupled*, and that therefore a *joint* algorithm has to be developed for this task. We propose both centralized and distributed algorithms that combine aspects of the above-mentioned methods.

For the implementation of *centralized* scheduling and streaming algorithms, the presence of cellular infrastructure can be exploited (e.g., see [35]), where the devices communicate directly but are under control of an existing cellular base station. The scheduling component of the proposed scheme is based on a link conflict graph (e.g., formed centrally by the base station) such that the links scheduled to be active simultaneously at any time-slot must form an independent set of such conflict graph. Then, the scheduling decisions at any given slot time can be formulated as a max-weight independent set (MWIS) problem, where the weights evolve dynamically as discussed later. The MWIS problem is known to be NP-hard; in this work, we make use of a message-passing algorithm proposed in [36] and [37] to efficiently obtain an approximate solution. Therefore, our *scheduling* component is designed based on this message-passing concept with D2D-related modifications, where the weights depend on both the lengths of the transmission queues and on the instantaneous link quality, which fluctuates due to the wireless channel fading.

The proposed *distributed* scheduling scheme is based on a modification of FlashLinQ, where the link priority is given by dynamically adjusted weights, which depend on the streaming process, and therefore create the connection between the transmission scheduling and streaming decisions for video quality.

Our *streaming* component is based on a stochastic network optimization framework, where the objective is to maximize the users' time-averaged video quality. Each video consists of a number of chunks. Each chunk can be requested at different quality levels,² such that higher quality corresponds to more bits per chunk to be delivered. Therefore, our algorithm dynamically controls the quality mode of each chunk to maximize total quality subject to all data being supportable over the network. Note that the streaming decisions impact the weights (for MWIS or modified FlashLinQ) of the scheduling.

Our scheme works with any source coding that provides sequential video chunks at different quality/rate levels, and can switch between levels at any chunk, depending on the control policy decisions. The proposed scheme in this paper also works when it is used in conjunction with multiple copies of the same video, encoded at different quality/rate levels, as it is currently done in Dynamic Adaptive Streaming over HTTP (DASH) and supported by common commercial services such as Netflix, Amazon Prime, and iTunes.

A. Related Work

D2D wireless networks are special cases of peer-to-peer (P2P) networks where the topology (who can communicate with whom) is determined by the radio propagation channel,

²For example, this can be obtained by storing multiple copies of the same video encoded at different rates, as in current video on-demand delivery systems such as Netflix or Amazon Prime, or by using scalable video coding and requesting more or fewer refinement layers [38]–[40].

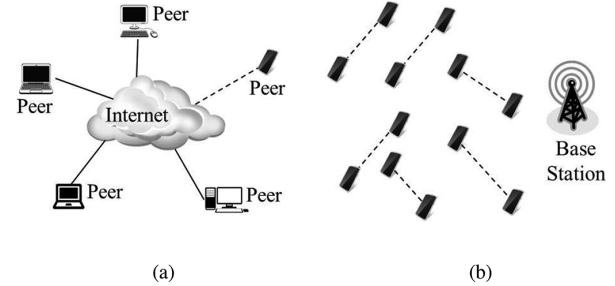


Fig. 1. Architectural difference between P2P networks and D2D networks: The solid and dashed lines represent wireline and wireless links, respectively. (a) P2P networks. (b) D2D networks.

which is a shared medium with interference. In particular, due to the distance-dependent path loss, in D2D networks a link may exist only if the two nodes are in the transmission range of each other (see Fig. 1). In contrast, in a general P2P network (see Fig. 1), a variety of topologies (e.g., tree, mesh) can be implemented as “overlay” over a common IP substrate (cloud), with the purpose of achieving better scalability, load balancing, and default tolerance than conventional client–server architectures [41]–[48]. A considerable amount of work in P2P networks has been devoted to game-theoretic frameworks, where nodes are given incentives to share their resources in order to avoid free-riding peers (e.g., [41], [42], and [49]). Source content replication in order to increase the peer resource sharing opportunities has also been widely investigated (e.g., [43], [45], and [50]).

Because D2D wireless networks are special cases of P2P wireless networks, some research results consider the terminologies of D2D and P2P as identical (e.g., see [51]–[53]). In [51], a scheme called “Microcast” is considered in order to implement efficient multicast by exploiting D2D short-range communication. In this system, a group of co-located devices collaborate by downloading different segment of the same file and then sharing such segments through network-coded local D2D links. This approach is relevant to the case of live-streaming, where all the users wish to stream the same video file at the same time. In contrast, in our work we consider on-demand video streaming where the users request different individual files. Hence, the Microcast architecture of [51] is completely ineffective in our case. In [52], the problem of file placement in a wireless P2P system with caching nodes is addressed. The considered system allows multihop, therefore file subpackets are placed such that the number of hops to reach their destination is proportional to the playback deadline delay in the streaming session. In our paper, we restrict to single-hop networks, and we address the queuing and transmission delay through the dynamic adaptation of the video quality level, instead of looking at the number of hops (which is equal to 1 for all streaming sessions in our case). The work in [53] treats a general base-station aided D2D network, where the D2D component is an underlay extension to the cellular network. The paper develops a general scheduling approach which is not targeted to video streaming and video quality adaptation.

B. Contributions

This work extends previous investigations of adaptive video streaming algorithms [54], [55]; we retain the notation of those

papers for streaming-related aspects. The algorithms in [54] and [55] are suited for adaptive stochastic video streaming by means of cache/helper dedicated nodes (e.g., access points connected to video servers through high-capacity wired links) and consider a bipartite network topology formed by helper and user nodes. Our work differs from the previous contributions mentioned before in the following aspects.

- The algorithms proposed in this paper extend the work in [54] and [55] in terms of: 1) new scheduling policies taking into account the D2D link conflict graph, and 2) a D2D-specific network model via multiple fixed source–destination pairs.
- In [54] and [55], it is assumed that each user can be served simultaneously by multiple infrastructure nodes over each scheduling slot. Instead, here we explicitly consider the constraint of the D2D link conflict graph, such that at each scheduling slot, a user can only be served by another (peered) user device if the corresponding link belongs to the scheduled independent set.
- The algorithms in [54] and [55] dynamically match source and destination pairs in every single unit time operation. However, the algorithms in this paper work on fixed source–destination pairs, as this is the relevant case for D2D communications. Let us note that FlashLinQ also considers this case [25], [26] where setting up D2D links dynamically and frequently during a given streaming session incurs in a too large protocol overhead.
- Since user devices are paired permanently over a whole streaming session in the present work, there is no handover delay, while per-slot dynamic association used in [54] and [55] gives rise to such delay. Such delay must be taken into account or made small through a special network architecture, e.g., single-channel single-IP implementation [56], [57].
- Differently from [41], [42], and [49], our work does not consider “strategic” users, which need incentives to cooperate. We assume that the nodes participating in the D2D network obey the rules of the designed protocol (incentives may exist at the level of the service provider, in order to make such D2D streaming service attractive from the customers viewpoint, but we are not concerned with this aspect here).
- Last, we evaluate the performance of our proposed algorithms by extensive simulations and compare them to a baseline scheme formed by FlashLinQ at the MAC layer and DASH at the application layer. In particular, we study the system performance in terms of total throughput, video quality (PSNR), and number of video streaming stall events³ at the receivers. According to our simulation results, the proposed algorithms provide sizeable performance gains for these quality measures.

The remainder of this paper is organized as follows: Section II gives preliminaries and background information. Section III explains the details of our proposed quality-aware streaming and scheduling algorithms both for the centralized

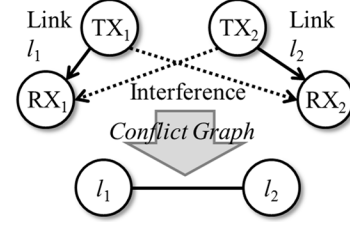


Fig. 2. Example of a conflict graph.

and the distributed cases. Section IV shows the simulation results compared to FlashLinQ variants. Section V concludes this paper.

II. PRELIMINARIES

This section defines our reference model including: 1) network model (see Section II-A); 2) link model (see Section II-B); and 3) wireless video streaming model (see Section II-C).

A. Reference Network Model: Macro View

Consider a network formed by a set \mathcal{L} of one-hop D2D links, indicated by $l_i \in \mathcal{L}$ [35], [58]. To schedule the D2D links, a conflict graph is constructed such that the set of vertices is \mathcal{L} (the links) and two vertices are connected by an edge if the corresponding links suffer from mutual interference above a certain desired threshold (refer to Fig. 2). The choice of the threshold will be discussed in Section IV-B.3. However, for the computation of the achievable rates of each link, we still need to take into account the residual interference, caused by the transmission from simultaneously scheduled links.

The conflict graph is described through its adjacency matrix, whose elements $\mathcal{E}_{(j,k)}$ between $l_j \in \mathcal{L}$ and $l_k \in \mathcal{L}$ are defined as follows:

$$\mathcal{E}_{(j,k)} = \begin{cases} 1, & \text{if } l_j \text{ interferes with } l_k \text{ where} \\ & l_j \in \mathcal{L}, l_k \in \mathcal{L}, \text{ and } j \neq k \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

In addition, the set of neighbor nodes of each node is defined as follows:

$$\mathcal{N}(i) \triangleq \{l_a | \mathcal{E}_{(i,a)} = 1 \text{ where } l_a \in \mathcal{L}\} \quad \forall l_i \in \mathcal{L}. \quad (2)$$

B. Reference Link Model: Micro View

As can be seen in Fig. 3, each D2D wireless link consists of one transmitter and its associated receiver. Each transmitter has a queue whose length evolves according to

$$Q_i(t+1) = \max[0, Q_i(t) - \mu_i(t)] + \lambda_i(t) \quad (3)$$

where $Q_i(t)$, $\mu_i(t)$, and $\lambda_i(t)$ stand for the queue backlog size at the transmitter of l_i , the number of bits leaving the queue of the transmitter of l_i , and the number of bits added to the queue of the transmitter of l_i , respectively. As shown in Fig. 3, the

³When the playback buffer does not contain the required video chunk at its due playback time such that playback has to stall and wait until such chunk is delivered.

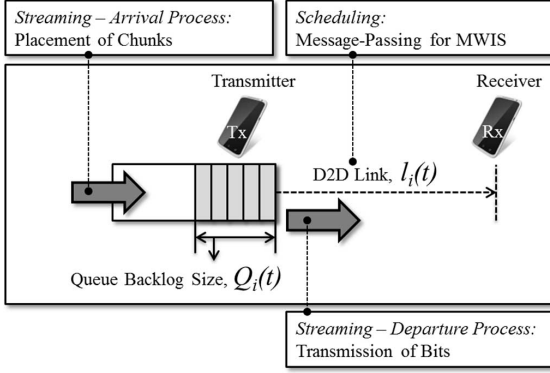


Fig. 3. D2D link model.

arrival process (bits added to the queue) is associated with the placement of chunks of the currently served video in each unit time $t \in \{0, 1, \dots\}$. When the D2D link is scheduled for transmission by the centralized or distributed controller, the queue has a departure process that depends on the channel state and on the link scheduling decisions. More details are provided in Section III.

C. Wireless Video Streaming System Model

The receiver of link l_i requests a video file f_i that is located in the cache of its associated transmitter. A video file is formed by a sequence of *chunks*, i.e., group of pictures (GOPs), which are encoded and decoded as standalone units. Chunks must be reproduced in sequence at the D2D receivers. The streaming thus consists of the transmission of sequential chunks from the transmitter to its associated receiver such that the playback buffer at each transmitter contains the required chunks at the beginning of each chunk playback time.

The time scale for the scheduling decision and departure process, i.e., t , is not equivalent to the chunk placement unit time τ , as can be seen in Fig. 4.

A chunk contains $\mathcal{N} = N_{\text{fpc}} \cdot N_{\text{ppf}}$ pixels, where N_{ppf} denotes the number of pixels per frame and N_{fpc} stands for the number of frames per chunk. Suppose that each chunk of each file f is encoded at a number of different quality modes $q \in M$ where $M = \{q_1 \dots q_M\}$. According to the variable bit-rate nature of video coding, the quality-rate profile may vary from chunk to chunk. We let $\mathbb{P}_f(q, \tau)$ and $\mathcal{N}\mathbb{B}_f(q, \tau)$ denote the video quality measure [e.g., peak-signal-to-noise-ratio (PSNR)]⁴ and the number of bits for file f at chunk time τ with quality mode q , respectively.

The chunk placement procedure consists of choosing the quality mode $q_i(\tau)$ of the chunks requested at chunk unit time τ by the scheduled D2D transmitters i . The choice $q_i(\tau)$ renders the point $(\mathbb{P}_{f_i}(q_i(\tau), \tau); \mathcal{N}\mathbb{B}_{f_i}(q_i(\tau), \tau))$ from the finite set of quality-rate tradeoff points $\{(\mathbb{P}_{f_i}(q, \tau), \mathcal{N}\mathbb{B}_{f_i}(q, \tau))\}_{q=q_1}^{q=q_M}$. The network controller: 1) chooses the quality mode $q_i(\tau)$ for chunk time τ for all requesting receivers i ; and 2) allocates the source coding rate (bit per pixel). The transmitter of link i

⁴There is a rich literature on video quality metrics, e.g., [33], [59], and [60]. Our framework works with any video quality measure, but for the sake of simplicity (and because details of quality measures are outside the scope of this paper), we use PSNR henceforth.

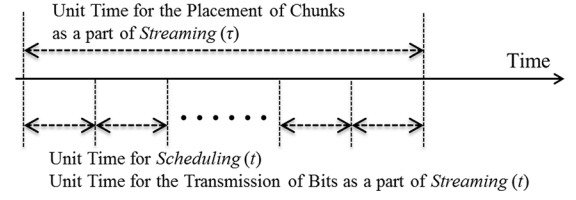


Fig. 4. Two differentiated unit time scales.

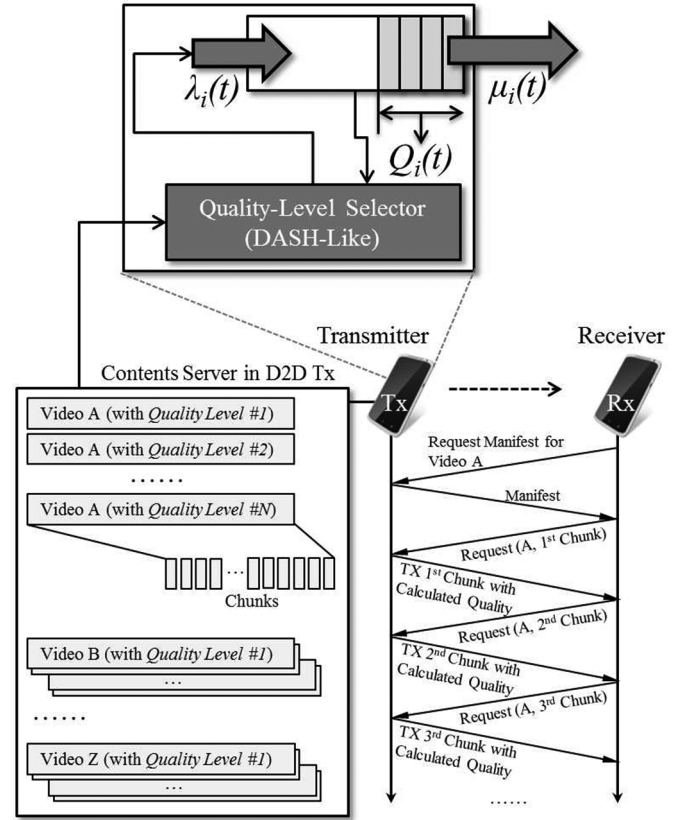


Fig. 5. Proposed DASH-like quality adaptive D2D system architecture. The D2D transmitter has its own contents server that contains various videos (with various quality levels), and each video consists of a sequence of chunks. For each D2D transmission, the transmitter can select suitable quality (calculated by the proposed algorithm in Section III-C) for each chunk transmission. The proposed algorithm in Section III-C can be represented as a *Quality-Level Selector (DASH-Like)* in this figure that determines which quality level/mode (among all possible quality levels) based on the transmitter queue backlog size.

places the corresponding $\mathcal{N}\mathbb{B}_{f_i}(q_i(\tau), \tau)$ bits in its transmission queue $Q_i(\tau)$, to be sent to the receiver whose length evolves according to.

Summarizing, as can be seen in Fig. 5, at each chunk time τ , the transmitter fetches chunks from its local cache in the order in which they are to be played back. The chunks are fetched at the quality mode $q_i(\tau)$ that is computed based on stochastic network optimization algorithms (details are in Section III-C.1). Then, the coded bits are packetized to be transmitted over the air interfaces. The enqueued packets will be transmitted depending on the departure process $\mu_i(t)$ and are dependent on channel states as well as interference from activated neighbor D2D links, i.e., signal-to-interference-plus-noise ratio (SINR).

III. QUALITY-AWARE STREAMING AND SCHEDULING FOR DEVICE-TO-DEVICE VIDEO DELIVERY

This section presents the basic design rationale of our proposed two quality-aware streaming and scheduling algorithms.

A. Design Rationale

The proposed centralized algorithm consists of two separable but interconnected parts, i.e., *centralized scheduling based on MWIS* (refer to Section III-B.1) and *quality-aware streaming* (refer to Section III-C). In order to (approximately) solve the MWIS problem in a computationally efficient manner, we resort to a message-passing approach. For the streaming decision, the operations to control the arrival and departure processes in the queue of each D2D transmitter are defined. The entire link model is illustrated in Fig. 3.

The proposed distributed algorithm improves FlashLinQ with the concepts of *distributed max-weight scheduling* before transmission (refer to Section III-B.2) and quality-aware streaming (refer to Section III-C).

B. Device-to-Device Scheduling

1) *Centralized Scheduling With MWIS Formulation*: For centralized D2D scheduling, the objective is to find the set of links (i.e., nodes of the conflict graph defined before) that maximize the sum of weights over all possible independent sets. This yields the MWIS problem

$$\max : \mathcal{F}(\mathcal{I}) \triangleq \sum_{\forall l_i \in \mathcal{L}} w_i \mathcal{I}_i, \quad (4)$$

$$\text{s.t. } \mathcal{I}_j + \mathcal{I}_k + \mathcal{E}_{j,k} \leq 2 \quad \forall l_j \in \mathcal{L}, \forall l_k \in \mathcal{L} \quad (5)$$

$$\mathcal{I}_i \in \{0, 1\} \quad \forall l_i \in \mathcal{L} \quad (6)$$

where \mathcal{I}_i is defined as

$$\mathcal{I}_i = \begin{cases} 1, & \text{if } l_i \text{ is scheduled where } l_i \in \mathcal{L} \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

The above formulation ensures that conflicting links are not scheduled simultaneously: If $\mathcal{E}_{j,k} = 0$ (no edge between l_j and l_k), then $\mathcal{I}_j + \mathcal{I}_k \leq 2$, i.e., both indicator functions can be equal to 1. In contrast, if $\mathcal{E}_{j,k} = 1$, $\mathcal{I}_j + \mathcal{I}_k \leq 1$, i.e., at most one of the two indicator functions can be equal to 1.

In (4), $w_i, \forall i \in \{1, \dots, |\mathcal{L}|\}$ is given by [61]

$$w_i \triangleq r_i(t) \cdot Q_i(t) \quad (8)$$

where $Q_i(t)$ is the queue backlog size at the transmitter of D2D link l_i , defined in (3), and $r_i(t)$ stands for the achievable rates of D2D link l_i . As formulated in (8), the definition of weights is associated with *instantaneous* queue backlog sizes that vary over time t . Therefore, this definition is different from the definition of the transmission throughput in P2P systems, which is by definition a time-averaged quantity.

The exact value of $r_i(t)$ of D2D link l_i cannot be obtained before a scheduling decision is made because it depends on the actual interference that all active transmitters (of the scheduled links) cause on the receiver of link l_i , which is known only when the scheduling decision is actually made.

- Input
 - w_i in (8) where $\forall l_i \in \mathcal{L}$
 - $\mathcal{E}_{(j,k)}$ in (1) where $\forall l_j \in \mathcal{L}, \forall l_k \in \mathcal{L}$
- Output
 - $\mathcal{F}(\mathcal{L}, \mathcal{E})$
 - \mathcal{L}^* // set of scheduled D2D links

Update Phase;

$n = 1$;

while $n \leq K$ **do**

// K : the number of message-passing iteration;

$m_{i \rightarrow j}^n = \max \left[0, w_i - \sum_{k \in \mathcal{N}(i)-j} m_{k \rightarrow i}^{n-1} \right], \forall j \in$

$\mathcal{N}(i)$;

i **sends** $m_{i \rightarrow j}^n$ **to all** $j \in \mathcal{N}(i)$;

$n++$;

end

Estimation Phase;

$$\mathcal{I}_i = \begin{cases} 1 & \text{if } \sum_{k \in \mathcal{N}(i)} m_{k \rightarrow i}^K < w_i \\ 0 & \text{otherwise} \end{cases};$$

If $\mathcal{I}_i = 1$ **then** $l_i \in \mathcal{L}^*$;

MWIS Computation Phase;

$$\mathcal{F}(\mathcal{L}, \mathcal{E}) = \sum_{\forall l_i \in \mathcal{L}} w_i \mathcal{I}_i;$$

Algorithm 1: MWIS-based scheduling with message-passing in each $l_i \in \mathcal{L}, \forall i \in \{1, \dots, |\mathcal{L}|\}$

To circumvent this problem, the decision in (4) is made on the basis of an estimated value of the link achievable rate, given by

$$r_i(t) = \log_2 \left(1 + \frac{\mathcal{P}_{s_i \rightarrow d_i}(t) \|h_{i \rightarrow i}\|^2}{\sigma^2 + \gamma} \right) \quad (9)$$

where $\mathcal{P}_{s_i \rightarrow d_i}(t)$ stands for the transmit power from s_i at unit time t , $h_{i \rightarrow i}$ stands for the (complex amplitude) channel gain from s_i to d_i , σ is the standard deviation of the (Gaussian) background noise, and γ stands for the interference thresholds, i.e., the maximum admissible interference level γ from a single interferer scheduled at the same time as the considered link. In FlashLinQ, γ is set to 9 dB [25], [26]. Equation (9) implies the assumption that the *aggregate* interference from other links is equal to the interference from a maximally strong single interferer. Since the overall interference levels tend to be dominated by the strongest interferer [62], this is a reasonable approximation.

After solving this MWIS problem, a set of active links is obtained, and the actual rates (including all the interference caused by the active transmitters on the link receivers) are used to update the transmission queues (see later). For solving the MWIS problem, various heuristic and approximation algorithms have been proposed due to the fact that MWIS is a well-known NP-hard problem. One of these methods is the computation with *message-passing* [36], [37], which we will use henceforth. The corresponding pseudo-code is presented in Algorithm 1.

2) *Distributed Max-Weight Scheduling*: For distributed scheduling, we improve FlashLinQ with the concept of max-weight scheduling.

FlashLinQ sorts the links according to a priority order externally determined (e.g., at random, or according to round robin) and considers the links in sequence, according to the priority order. A link is declared active if it does not create significant interference to the already active links with higher priority and if its own achieved rate, considering the interference from the already active links with higher priority, is large enough. The decision is based on measurements of channel strengths in both directions [25], [26]. In the original FlashLinQ, priorities are randomized over time to provide a basic level of fairness. With the concept of max-weight scheduling, we set the priorities of D2D links instead as follows:

$$U_i \triangleq \frac{1}{r_i(t) \cdot Q_i(t)}. \quad (10)$$

C. Streaming With Quality-Aware Stochastic Control

The streaming consists of two parts, i.e.: 1) placement of chunks (i.e., arrival process of the queue), and 2) transmission of bits (i.e., departure process of the queue). Notice that the streaming method investigated in this section is used for both centralized and distributed algorithms.

1) *Arrival Process (Placement of Chunks)*: In each chunk time-slot $\tau \in \{0, 1, \dots\}$, the transmitter of each link places a chunk into its transmission queue.

In order to dynamically and adaptively select the quality level of the chunks, we consider the following stochastic optimization approach that aims at maximizing the total average video quality of the users. Let $\mathbb{P}(t) = \sum_{l_i \in \mathcal{L}} \mathbb{P}_{f_i}(q_i(t), t)$, and

$$\mathbb{P}_{f_i}(q_i(t), t) = \begin{cases} \mathbb{P}_{f_i}(q_i(t), t), & \tau \bmod t = 0 \\ 0, & \tau \bmod t \neq 0. \end{cases} \quad (11)$$

Then, the proposed stochastic optimization problem is given by

$$\max \quad \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{t^*=0}^{t-1} \mathbb{E}[\mathbb{P}(t^*)] \quad (12)$$

$$\text{subject to} \quad \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{t^*=0}^{t-1} \mathbb{E}[Q_i(t^*)] < \infty \quad \forall l_i \in \mathcal{L} \quad (13)$$

where (13) means all given queues should fulfill mean rate stability. Let $\Theta(t)$ denote the column vector of all scheduled queues at time t , and define the quadratic Lyapunov function

$$L(t) = \frac{1}{2} \Theta^T(t) \Theta(t) = \frac{1}{2} \sum_{i \in \mathcal{L}} |Q_i(t)|^2 \quad (14)$$

where $\Theta^T(t)$ stands for the transpose of $\Theta(t)$. Then, let $\Delta(t)$ be a conditional quadratic Lyapunov function that can be formulated as $\mathbb{E}[L(t+1)|\Theta(t)] - L(t)$, i.e., the drift on slot t . The *drift-plus-penalty* (DPP) policy is designed to solve the given optimization formulation by observing only the current queue backlog sizes $\Theta(t)$ and choose q (i.e., quality mode) to maximize a bound on

$$\mathbb{P}(t) - \alpha \Delta(t) \quad (15)$$

where α is a positive constant control parameter of the DPP policy that affects the quality-delay tradeoffs [61].

Now, the quality control decision involves choosing $q_i(t)$, the quality mode for all scheduled receivers at chunk time t^5 . This choice is made as

$$\arg \max_{q_i(t) \in M} [\mathbb{P}_{f_i}(q_i(t), t) - \alpha \{\mathcal{N}\mathbb{B}_{f_i}(q_i(t), t)\} Q_i(t)]. \quad (16)$$

Since the placement of chunks constitutes the arrival process of the queue, $\lambda_i(t)$ can be denoted as follows when the optimal q is determined in each $l_i \in \mathcal{L}$:

$$\lambda_i(t) = \begin{cases} \mathcal{N}\mathbb{B}_{f_i}(q_i(t), t), & \tau \bmod t = 0 \\ 0, & \tau \bmod t \neq 0. \end{cases} \quad (17)$$

2) *Departure Process (Transmission of Bits)*: Once a set of active links \mathcal{L}^* is determined as described in Section III-B, the transmitters of the scheduled links can transmit bits up to the amount of the achievable rates actually supported by the link at time t , i.e., $\mu_i(t) = r_i(t)$, $\forall l_i \in \mathcal{L}^*$.

According to Shannon's capacity equation, i.e., $\mu_i(t)$ in (8) can be computed as follows [63]:

$$\mu_i(t) = \mathcal{B} \cdot \log_2 \left(1 + \frac{\mathcal{P}_{s_i \rightarrow d_i}(t) \|h_{i \rightarrow i}\|^2}{\sigma^2 + \sum_{j \neq i} \mathcal{P}_{s_j \rightarrow d_i}(t) \|h_{j \rightarrow i}\|^2} \right) \quad (18)$$

where $\forall l_i \in \mathcal{L}^*, \forall l_j \in \mathcal{L}^*, i \neq j$, $\mathcal{P}_{s_a \rightarrow d_b}(t)$ stands for the power transmitted by s_a intended for d_b , and $h_{j \rightarrow i}$ stands for the channel gain from the transmitter of link j to the receiver of link i where $\forall a, \forall b \in \{1, \dots, |\mathcal{L}|\}$ at time t , \mathcal{B} stands for the channel bandwidth of the system.

While here we have used the SINR-based capacity equation in (18) to evaluate achievable rates, any suitable function of SINR can be included in our algorithms, for example, if the physical layer (PHY) of the D2D system includes a family of modulation and coding scheme (MCS), each of which has a certain operational range of SINR and a given rate, we can substitute such a piecewise constant function into our scheme and get meaningful results that explicitly include the properties of the MCS set (e.g., the MCS modes of 802.11-based standards, or 3GPP LTE).

IV. SIMULATION STUDY

The performance of our proposed joint scheduling and streaming algorithms is simulated and evaluated in this section. The basic simulation settings are presented in Section IV-A, and the simulation results are presented and analyzed in Section IV-B.

A. Simulation Settings

1) *Video Traces*: For the simulation study, we use four different types of video traces. Each video consists of 14 400 chunks where the playback time of each chunk is 0.5 s. Thus, the overall playback time of each video trace is 2 h, corresponding to a typical movie playback time. The video

⁵This quality mode $q_i(t)$ selection decision also determines $\Phi(t)$ the column vector of the number of source-coded bits $\mathcal{N}\mathbb{B}_{f_i}(q_i(t), t)$ with selected quality mode $q_i(t)$ that each receiver i must download from its transmitter for the chunk requested at time t .

TABLE I
VIDEO TRACE INFORMATION

	Basic Information (Names of Test Sequences)	Resolution	Average Bitrates (Full Video Stream with All Layers)
Video Trace 1	highway	352 × 288 Pixels	631 Kbps (8 different quality levels available)
Video Trace 2	city, crew, harbour, train	704 × 576 Pixels	3908 Kbps (4 different quality levels available)
Video Trace 3	parkrun, stockholm	1024 × 576 Pixels	6679 Kbps (4 different quality levels available)
Video Trace 4	bridge_close, bridge_far	352 × 288 Pixels	556 Kbps (8 different quality levels available)

sequences are standard Moving Picture Experts Group (MPEG) test sequences, commonly used in the literature. The original video sequences consist of 200 chunks. To create one 2-h video, we concatenated the same sequence 72 times. Details of the traces are summarized in Table I.

The quality of each chunk can be numerically represented by the PSNR, i.e., $\mathbb{P}_{f_i}(q_i(t), t)$.

We note that the video streams are not synchronized between the D2D links, i.e., starting times for the different links (streams) are independent of each other.

2) *Baseline Terms of Comparison:* To show the effectiveness of our proposed centralized or distributed quality-aware streaming and scheduling algorithms, their performances are evaluated and compared to the performances of FlashLinQ variants.

- *FlashLinQ:* This is the standard FlashLinQ algorithm with random U_i priority selection in each D2D link, such that each link has at least $\frac{1}{N}$ probability of being scheduled, where N stands for the number of D2D links. Since FlashLinQ does not consider video quality at all, we fix the quality level as 2 in video trace 1 (among given 4 levels), 4 in video trace 2 (among given 8 levels), 4 in video trace 3 (among given 8 levels), and 2 in video trace 4 (among given 4 levels). We choose a *medium* quality level since a selected high quality level leads to high PSNR but might negatively impact queue stability.
- *FlashLinQ-P:* This variant of FlashLinQ uses prioritized U_i selection. The U_i are computed as

$$U_i = \frac{1}{r_i(t) \cdot Q_i(t)} \quad (19)$$

corresponding to the max-weight scheduling concept. Also *FlashLinQ-P* does not consider video-quality related aspects, and we again fix the qualities of the streams to the same values as above.

- *FlashLinQ-Q:* This variant of FlashLinQ uses the video streaming quality decisions as in Section III-C, but the link scheduling is standard FlashLinQ, i.e., the priorities U_i are chosen at random.

3) *Simulation Topology Construction:* The considering simulation topology consists of uniformly random deployed 10 D2D transmitter and receiver pairs in a $600 \times 600\text{-m}^2$ square layout. The path loss is computed according to the Winner II model (indoor in 2.4 GHz D2D communications) [24]

where f_c^{GHz} is the carrier frequency in a GHz scale. a_1 includes the path loss exponent, and its value is 18.7 dBm in LOS and 36.8 dBm in NLOS. a_2 is the intercept, which is 46.8 dBm in LOS and 43.8 dBm in NLOS. a_3 describes the path-loss frequency dependence, and it is set to 20 in both LOS and NLOS. X_σ is the shadowing assumed to be a normal distribution (in dB) with mean 0 and standard deviation σ , where $\sigma = 3$ dB in LOS and $\sigma = 6$ dB in NLOS. Notice that we assume that no communication is possible for a distance longer than 100 m.

With the given topology, we simulated transmission for 24 h, assuming each D2D link streams 10 video traces. Then, 10 simulations will be operated with different random geometries.

The connectivity of conflict graphs changes according to the setting of the interference threshold. If the received powers from nearby D2D transmitters j to current D2D receiver i are less than the interference threshold, they will be considered as noise. Otherwise, they will be considered as interference. Thus, low γ increases the number of edges in the corresponding conflict graph. Consequently, a relatively small number of D2D links will be scheduled; that may also reduce the sum rate. On the other hand, high γ decreases the number of edges in the given corresponding conflict graph. Thus, a relatively large number of D2D links can be scheduled, however it will reduce the signal-to-interference-plus-noise ratio, and thus rate, for each link. Therefore, we need to consider various interference threshold settings in the simulation studies.

4) *Performance Analysis Metrics:* Whenever a receiver finishes playing back chunk i , all bits of chunk $i + 1$ should have arrived. Otherwise, a pause (*stall*) occurs in the playback; see Fig. 6. Obviously, the stall events introduce user dissatisfaction. Typically, three or more stalls during one movie would be judged to be unacceptable quality. Thus, the number of stall events can be an important index for measuring user satisfaction of video streaming.

To avoid stall events, *pre-buffering* is frequently used. As shown in Fig. 7, a number of chunks are transmitted before playback at the receiver starts; this introduces a viewing delay for the user, but reduces the number of stall events. Obviously a very long pre-buffering time leads to user dissatisfaction as well, and if we set extremely large pre-buffering time, there is no difference between streaming and downloading/transmission. Therefore, there exists a tradeoff, and defining an appropriate pre-buffering time is required. In addition, once a stall occurs, the receiver buffers again for the same amount of pre-buffering time as for startup.⁶

$$PL(d) = a_1 \log_{10}(d) + a_2 + a_3 \log_{10} \left(\frac{f_c^{\text{GHz}}}{5} \right) + X_\sigma \quad (20)$$

⁶Notice that the re-buffering time could also be shorter than the initial buffering time, i.e., they do not need to be the same, but we set it as equal in order to reduce the number of variables in the simulation.

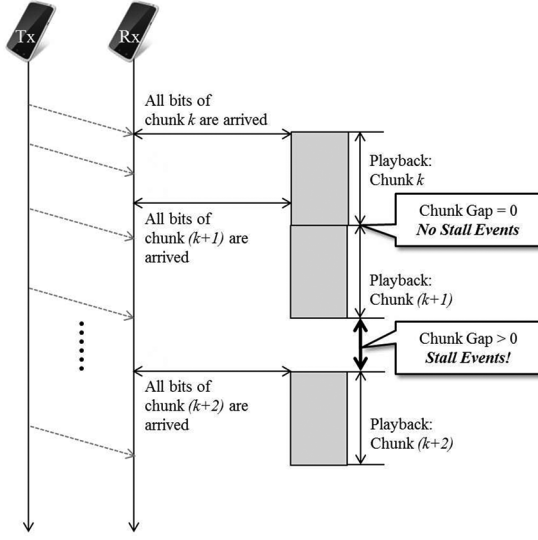


Fig. 6. Example of stall events: If all bits of next chunk are arrived at the playback buffer of D2D receiver, there is no stall event since the receiver can immediately play the next chunk when the playing of current chunk is completed. Otherwise, the stall event will occur when the playing of current chunk is completed.

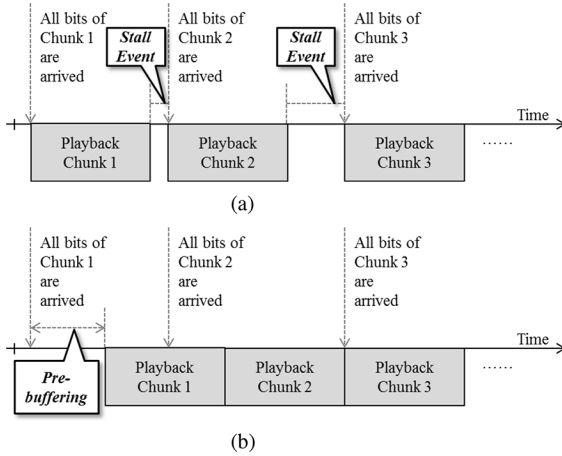


Fig. 7. Pre-buffering: With the definition of pre-buffering, the number of stall events can be reduced, i.e., user satisfaction can be increased. (a) No pre-buffering: If there is no pre-buffering, stall events may occur between chunks. (b) Pre-buffering: By setting certain amounts of pre-buffering time, we can reduce the number of stall events.

B. Simulation Results

With the given two performance metrics, we evaluate the performance of our proposed algorithms and three various FlashLinQ variants as a function of the following parameters:

- various pre-buffering time settings (see Section IV-B.1);
- various α , i.e., quality-delay tradeoffs (see Section IV-B.2);
- various interference thresholds γ (see Section IV-B.3);
- average quality versus the expected number of stall events (see Section IV-B.4).

Notice that our proposed centralized algorithm is denoted as mpMWIS-QP (i.e., *message-passing for MWIS* formulation with *Quality-awareness* and *max-weight Prioritization*); and

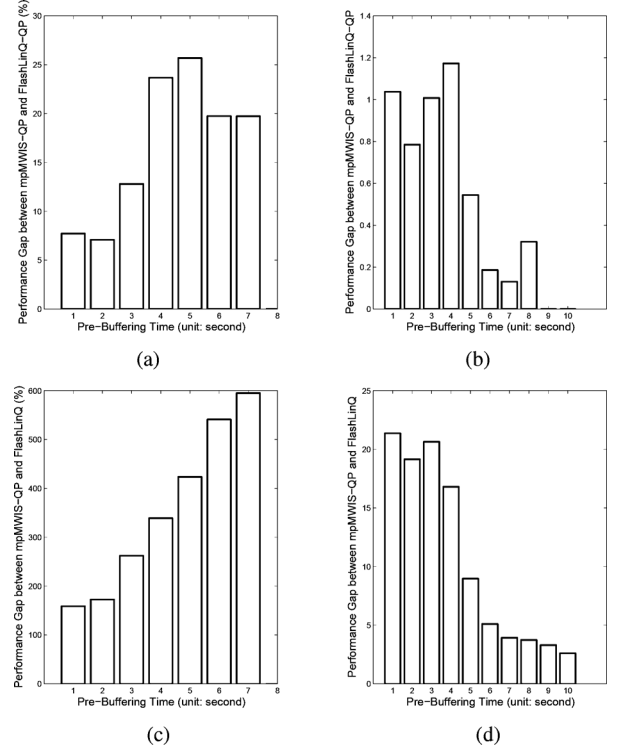


Fig. 8. Performance comparison between mpMWIS-QP and FlashLinQ-QP/FlashLinQ in terms of the expected number of stall events. (a) Effectiveness of mpMWIS-QP compared to FlashLinQ-QP in terms of \mathcal{T}_s defined in (21). (b) Effectiveness of mpMWIS-QP compared to FlashLinQ-QP in terms of \mathcal{M}_s defined in (22). (c) Effectiveness of mpMWIS-QP compared to FlashLinQ in terms of \mathcal{T}_s defined in (21). (d) Effectiveness of mpMWIS-QP compared to FlashLinQ in terms of \mathcal{M}_s defined in (22).

also our proposed distributed algorithm is denoted as FlashLinQ-QP (i.e., improved *FlashLinQ* with *Quality-awareness* and *max-weight Prioritization*).

1) *Effects of the Pre-Buffering Time*: We vary the pre-buffering time from 1 to 10 s with a step size of 1 s. In addition, α and the interference threshold γ are set to 2 and 5 dB, respectively. The resulting expected stall probability is presented in Table II and Fig. 8. For quantitative comparison, two indices, i.e., \mathcal{T}_s and \mathcal{M}_s , are defined respectively as

$$\mathcal{T}_s = \frac{\mathbb{E}[N_s] \text{ of FlashLinQ Variants} - \langle \mathbb{E}[N_s] \rangle}{\langle \mathbb{E}[N_s] \rangle} \quad (21)$$

$$\mathcal{M}_s = \mathbb{E}[N_s] \text{ of FlashLinQ Variants} - \langle \mathbb{E}[N_s] \rangle \quad (22)$$

where $\mathbb{E}[N_s]$ stands for the expected number of stall events and $\langle \mathbb{E}[N_s] \rangle$ denotes the $\mathbb{E}[N_s]$ of mpMWIS-QP.

As anticipated, the expected number of stall events reduces as the pre-buffering time increases. If there is no pre-buffering time, mpMWIS-QP has 13.4 stall events on average, whereas FlashLinQ has 34.8. For 8 s pre-buffering time, mpMWIS-QP has no stall events; FlashLinQ-QP has 0.3; this performance is the best among the given three FlashLinQ variants. Pure FlashLinQ shows the lowest performance. As shown in Fig. 8(c), it has 6 times more stall events when the pre-buffering time is 7 s.

We furthermore see that the performance advantage of mpMWIS-QP versus FlashLinQ stems from a variety of

TABLE II
EXPECTED NUMBER OF STALL EVENTS, I.E., $\mathbb{E}[N_s]$, IN EACH VIDEO STREAMING IN EACH D2D LINK WITH VARIOUS PRE-BUFFERING TIMES

Pre-buffering Time	mpMWIS-QP: $\mathbb{E}[N_s]$	FlashLinQ: $\mathbb{E}[N_s]$	FlashLinQ-P: $\mathbb{E}[N_s]$	FlashLinQ-Q: $\mathbb{E}[N_s]$	FlashLinQ-QP: $\mathbb{E}[N_s]$
1 second	13.4	34.8	26.0	26.6	14.5
2 second	11.1	30.2	20.9	22.9	11.9
3 second	7.9	28.5	15.5	18.3	8.9
4 second	5.0	21.8	10.2	11.3	6.1
5 second	2.1	11.1	4.4	5.7	2.7
6 second	0.9	6.0	2.2	3.0	1.1
7 second	0.7	4.6	1.6	2.4	0.8
8 second	0 [No Stalls]	3.7	1.2	1.9	0.3
9 second	0 [No Stalls]	3.3	0.6	1.1	0 [No Stalls]
10 second	0 [No Stalls]	2.6	0 [No Stalls]	0.3	0 [No Stalls]

TABLE III
EXPECTED NUMBER OF STALL EVENTS, I.E., $\mathbb{E}[N_s]$, IN EACH VIDEO STREAMING IN EACH D2D LINK WITH VARIOUS α

α	mpMWIS-QP: $\mathbb{E}[N_s]$	FlashLinQ: $\mathbb{E}[N_s]$	FlashLinQ-P: $\mathbb{E}[N_s]$	FlashLinQ-Q: $\mathbb{E}[N_s]$	FlashLinQ-QP: $\mathbb{E}[N_s]$
8	0 [No Stalls]	0.8	0.3	0.5	0 [No Stalls]
4	0 [No Stalls]	1.6	1.1	0.9	0.1
2	0 [No Stalls]	3.7	1.2	1.9	0.3
0.1	1.2	7.0	3.7	4.2	1.6
0.05	3.0	15.3	8.4	9.0	4.7

TABLE IV
EXPECTED NUMBER OF STALL EVENTS, I.E., $\mathbb{E}[N_s]$, IN EACH VIDEO STREAMING IN EACH D2D LINK WITH VARIOUS INTERFERENCE THRESHOLDS γ

γ	mpMWIS-QP: $\mathbb{E}[N_s]$	FlashLinQ: $\mathbb{E}[N_s]$	FlashLinQ-P: $\mathbb{E}[N_s]$	FlashLinQ-Q: $\mathbb{E}[N_s]$	FlashLinQ-QP: $\mathbb{E}[N_s]$
0 dB	5.4	26.9	20.8	22.0	17.9
5 dB	0 [No Stalls]	3.7	1.2	1.9	0.3
13 dB	3.1	23.8	13.0	16.2	10.2

factors: max-weight scheduling, incorporation of the interconnection between scheduling and streaming, and centralized control. We see that only incorporating max-weight scheduling (i.e., going from FlashLinQ to FlashLinQ-P) provides a significant advantage, while only incorporating video quality without max-weight scheduling (i.e., going from FlashLinQ to FlashLinQ-P) gives slightly lower improvement. The advantage of centralized scheduling over distributed scheduling (mpMWIS-QP) is very small, which is an important insight for actual deployment.

2) *Effects of the Parameter α* : α stands for a quality-delay tradeoff constant as formulated in (15). If α is small, quality awareness takes higher priority. On the other hand, larger α considers queue stability with higher priority, so that lower probability of stalls can be anticipated. The simulation results in this section investigate results when α takes on the values 0.05, 0.1, 2, 4, and 8. Notice that our considered pre-buffering time is 8 s, which is an optimum value for mpMWIS-QP; and the interference threshold γ is set to 5 dB. A further discussion of video quality versus stall events will be given in Section IV-B.4.

The simulation results are summarized in Table III. If α is 2, 4, or 8, there are no stall events in mpMWIS-QP. Pure FlashLinQ has between 3.7 and 0.8 stall events for those values. FlashLinQ-QP will have no stalls when $\alpha = 8$. This performance is lower than the performance of mpMWIS-QP, however the performance of FlashLinQ-QP is the best among the given FlashLinQ variants.

3) *Choice of the Interference Threshold γ* : As discussed in Section IV-A.3, the interference threshold trades off the number

TABLE V
PSNR TABLE OF GIVEN FOUR VIDEO TRACES

	minimum PSNR	Maximum PSNR
Video Trace 1	29.4835 dB	37.8063 dB
Video Trace 2	25.7136 dB	36.2584 dB
Video Trace 3	24.3273 dB	35.0470 dB
Video Trace 4	28.8283 dB	37.1691 dB

of active links with the rate per link that can be obtained. Thus, finding an appropriate interference threshold is important for optimizing performance. In our given network geometry, our minimum and maximum interference thresholds are 0 and 13 dB, respectively. With $\gamma = 0$ dB, our geometry is extremely densely connected in its corresponding conflict graph, i.e., only one D2D link will be scheduled in each unit time. On the other hand, our geometry has no edges in its corresponding conflict graph when $\gamma = 13$ dB, i.e., all D2D links will be scheduled and will generate interference all together in each unit time. We additionally performed the simulation with $\gamma = 5$ dB.

We again set the pre-buffering time to 8 s, and set $\alpha = 2$. Results with the mentioned three interference thresholds are listed in Table IV. For all algorithms, performance with $\gamma = 5$ dB is the best; the performance with $\gamma = 0$ dB is lower than the performance with $\gamma = 13$ dB.

4) *Average Quality Versus Expected Number of Stall Events*: This section presents average quality values depending on the expected number of stall events for mpMWIS-QP, FlashLinQ-Q, and FlashLinQ-QP. The other two FlashLinQ variants,

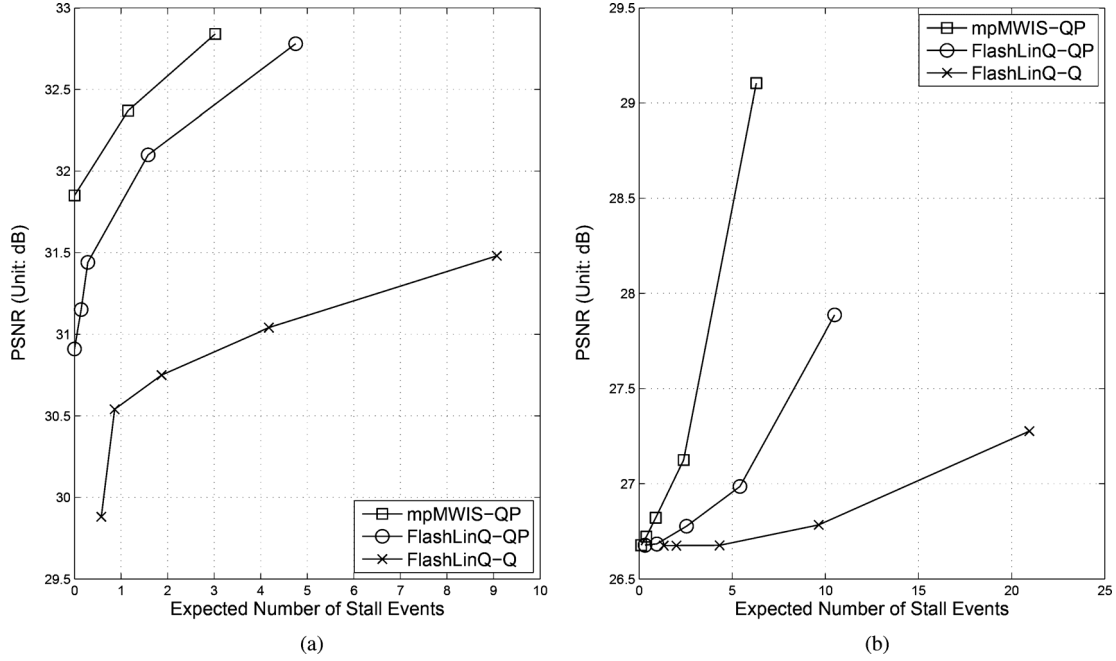


Fig. 9. Average quality versus the expected number of stall events. (a) System bandwidth: 2 MHz (plotted with the data in Table VI). (b) System bandwidth: 1 MHz (plotted with the data in Table VII).

TABLE VI
EXPECTED NUMBER OF STALL EVENTS $\mathbb{E}[N_s]$ VERSUS AVERAGE QUALITY (AVERAGE PSNR) WHEN THE SYSTEM BANDWIDTH IS 2 MHz

α	mpMWIS-QP	FlashLinQ-QP	FlashLinQ-Q
0.05	$\mathbb{E}[N_s]: 3.02$, Average PSNR: 32.84	$\mathbb{E}[N_s]: 4.75$, Average PSNR: 32.78	$\mathbb{E}[N_s]: 9.07$, Average PSNR: 31.48
0.1	$\mathbb{E}[N_s]: 1.15$, Average PSNR: 32.37	$\mathbb{E}[N_s]: 1.58$, Average PSNR: 32.10	$\mathbb{E}[N_s]: 4.18$, Average PSNR: 31.04
2	$\mathbb{E}[N_s]: 0.00$ [No Stalls], Average PSNR: 31.85	$\mathbb{E}[N_s]: 0.29$, Average PSNR: 31.44	$\mathbb{E}[N_s]: 1.87$, Average PSNR: 30.75
4	-	$\mathbb{E}[N_s]: 0.14$, Average PSNR: 31.15	$\mathbb{E}[N_s]: 0.86$, Average PSNR: 30.54
8	-	$\mathbb{E}[N_s]: 0.00$ [No Stalls], Average PSNR: 30.91	$\mathbb{E}[N_s]: 0.58$, Average PSNR: 29.88

TABLE VII
EXPECTED NUMBER OF STALL EVENTS $\mathbb{E}[N_s]$ VERSUS AVERAGE QUALITY (AVERAGE PSNR) WHEN THE SYSTEM BANDWIDTH IS 1 MHz

α	mpMWIS-QP	FlashLinQ-QP	FlashLinQ-Q
0.05	$\mathbb{E}[N_s]: 6.29$, Average PSNR: 29.10	$\mathbb{E}[N_s]: 10.50$, Average PSNR: 27.89	$\mathbb{E}[N_s]: 20.95$, Average PSNR: 27.28
0.1	$\mathbb{E}[N_s]: 2.39$, Average PSNR: 27.12	$\mathbb{E}[N_s]: 5.41$, Average PSNR: 26.99	$\mathbb{E}[N_s]: 9.65$, Average PSNR: 26.78
2	$\mathbb{E}[N_s]: 0.89$, Average PSNR: 26.82	$\mathbb{E}[N_s]: 2.55$, Average PSNR: 26.78	$\mathbb{E}[N_s]: 4.32$, Average PSNR: 26.68
4	$\mathbb{E}[N_s]: 0.39$, Average PSNR: 26.72	$\mathbb{E}[N_s]: 0.95$, Average PSNR: 26.68	$\mathbb{E}[N_s]: 2.00$, Average PSNR: 26.67
8	$\mathbb{E}[N_s]: 0.16$, Average PSNR: 26.68	$\mathbb{E}[N_s]: 0.32$, Average PSNR: 26.68	$\mathbb{E}[N_s]: 1.32$, Average PSNR: 26.67

i.e., FlashLinQ and FlashLinQ-P, are not considered in this simulation since they statically select their quality mode. Pre-buffering time is 8 s, and the interference threshold is 5 dB. To numerically represent video quality, PSNR is used; the minimum and maximum PSNR values in each video trace are listed in Table V.

We simulate the algorithms with $\alpha = \{0.05, 0.1, 2, 4, 8\}$, and compute the expected numbers of stall events and the average PSNR values. Results are shown for a system bandwidth of 2 MHz in Table VI. In mpMWIS-QP, there are no stall events if $\alpha \geq 2$. However, higher α leads to the degradation of average PSNR to guarantee more stability on the D2D transmitter queue. If there are no stall events, there is no need to improve the stability of the D2D transmitter queue. Thus, simulation results of mpMWIS-QP where $\alpha = 4$ and $\alpha = 8$ are not shown. In addition, FlashLinQ-QP has no stall events when $\alpha = 8$.

Fig. 9(a) shows that our mpMWIS-QP provides the highest PSNR for a given stall probability; the performance of FlashLinQ-QP is approximately 0.4 dB lower than the performance of mpMWIS-QP. However, FlashLinQ-Q shows around 1.6 dB lower PSNR compared to the PSNR of mpMWIS-QP.

Results for a system bandwidth of 1 MHz are in Table VII. Due to the lower bandwidth, all the three evaluated algorithms have stall events. Similar to the cases of 2 MHz system bandwidth, higher α leads to the degradation of average PSNR to guarantee more stability on the D2D transmitter queue. Fig. 9(b) shows that mpMWIS-QP again provides the highest PSNR, and the PSNR of FlashLinQ-QP is approximately 1.5 dB lower than the performance of mpMWIS-QP when the expected number of stall events is near 5. However, FlashLinQ-Q shows around 2.3 dB lower PSNR compared to the PSNR of mpMWIS-QP. According to Table VII, the lowest PSNR is

near 26.7, which is close to the PSNR of the lowest-quality mode available.

As observed in Fig. 9, a higher expected number of stall events is associated with a higher PSNR (i.e., video quality). For guaranteeing more video quality, lower α is used for putting more weights on PSNR and lower weights on queue stability. Therefore, there are more possibilities to increase the queue backlog sizes at D2D transmitters, i.e., this leads to higher expected number of stall events at D2D receivers.

V. CONCLUSION AND FUTURE WORK

This paper proposed centralized or distributed quality-aware streaming and scheduling algorithms that can be used for device-to-device video delivery applications. In terms of scheduling, we have considered both a centralized and a distributed approach. For centralized scheduling, a message-passing-based algorithm is used to obtain the solutions from a maximum independent set problem formulation. For distributed scheduling, we improved the FlashLinQ D2D scheduler by introducing a max-weight priority across the links. In terms of streaming, a quality-aware stochastic chunk selection algorithm is introduced that works based on the queue backlog sizes in each D2D transmitter queue. The stochastic algorithm in the streaming part controls the quality of each video chunk to maximize the qualities of streamed video subject to queue rate stability. We can draw several important conclusions from the simulations: 1) it is essential to use a transmission scheme that accounts for the interrelationship between scheduling and quality selection; 2) a good distributed scheme performs only marginally worse than our centralized scheme; and 3) we can trade off average video quality with probability of stalls. These results give important insight in the deployment of D2D-based video streaming.

As future research directions, the following two research problems are considerable.

- The proposed streaming algorithm is with stochastic network optimization for the tradeoff between *video quality* and *transmission queue stability*. Therefore, energy awareness is out of scope in this research, which can be a burden in power-hungry mobile devices. Thus, this factor can be additionally considerable in our future research.
- The real-world prototyping of current algorithms is considerable. As presented in [64], the earlier stage of the proposed algorithms was implemented with Android mobile platforms, and now the next-step implementation is under discussion for further improvement.

REFERENCES

- [1] Cisco, "Cisco visual networking index: Forecast and methodology 2012–2017," Cisco White Paper, 2013.
- [2] N. Golrezai, A. F. Molisch, A. G. Dimakis, and G. Caire, "Femto-caching and device-to-device collaboration: A new architecture for wireless video distribution," *IEEE Commun. Mag.*, vol. 51, no. 4, pp. 142–149, Apr. 2013.
- [3] S. Singh, J. G. Andrews, and G. de Veciana, "Interference shaping for improved quality of experience for real-time video streaming," *IEEE J. Sel. Areas Commun.*, vol. 30, no. 7, pp. 1259–1269, Aug. 2012.
- [4] A. Abdel Khalek, C. Caramanis, and R. W. Heath, Jr., "Video-aware MIMO precoding with packet prioritization and unequal modulation," in *Proc. EUSIPCO*, Bucharest, Romania, Aug. 2012, pp. 1905–1909.
- [5] S. Singh, O. Oyman, A. Papathanassiou, D. Chatterjee, and J. G. Andrews, "Video capacity and QoE enhancements over LTE," in *Proc. IEEE ICC ViOpt*, Ottawa, ON, Canada, Jun. 2012, pp. 7071–7076.
- [6] V. Joseph and G. de Veciana, "Jointly optimizing multi-user rate adaptation for video transport over wireless systems: Mean-fairness-variability tradeoffs," in *Proc. IEEE INFOCOM*, Orlando, FL, USA, Mar. 2012, pp. 567–575.
- [7] A. A. Khalek, C. Caramanis, and R. W. Heath, Jr., "Joint source-channel adaptation for perceptually optimized scalable video transmission," in *Proc. IEEE GLOBECOM*, Houston, TX, USA, Dec. 2011, pp. 1–5.
- [8] Z. Lu and G. de Veciana, "Opportunistic transport for stored video delivery over wireless networks: Optimal anticipative and causal approximations," in *Proc. 49th Allerton Conf. Commun., Control, Comput.*, Chicago, IL, USA, Oct. 2011, pp. 143–150.
- [9] C. Chen, R. W. Heath, Jr., A. C. Bovik, and G. de Veciana, "Adaptive policies for real-time video transmission: A Markov decision process framework," in *Proc. IEEE ICIP*, Brussels, Belgium, Sep. 2011, pp. 2249–2252.
- [10] L. Toni, P. C. Cosman, and L. B. Milstein, "Channel coding optimization based on slice visibility for transmission of compressed video over OFDM channels," *IEEE J. Sel. Areas Commun.*, vol. 30, no. 7, pp. 1172–1183, Aug. 2012.
- [11] L. Toni, P. C. Cosman, and L. B. Milstein, "Subcarrier mapping based on slice visibility for video transmission over OFDM channels," in *Proc. Asilomar Conf. Signals, Syst. Comput.*, Pacific Grove, CA, USA, Nov. 2012, pp. 920–924.
- [12] H. Ahlehagh and S. Dey, "Hierarchical video caching in wireless cloud: Approaches and algorithms," in *Proc. IEEE ICC ViOpt*, Ottawa, ON, Canada, Jun. 2012, pp. 7082–7087.
- [13] H. Ahlehagh and S. Dey, "Video caching in radio access network: Impact on delay and capacity," in *Proc. IEEE WCNC*, Paris, France, Apr. 2012, pp. 2276–2281.
- [14] S. Dey, "Cloud mobile media: Opportunities, challenges, and directions," in *Proc. IEEE ICNC*, Maui, HI, USA, Jan. 2012, pp. 929–933.
- [15] D. Wang, P. C. Cosman, and L. B. Milstein, "Cross layer resource allocation design for uplink video OFDMA wireless systems," in *Proc. IEEE GLOBECOM*, Houston, TX, USA, Dec. 2011, pp. 1–6.
- [16] L. Toni, P. C. Cosman, and L. Milstein, "Unequal error protection based on slice visibility for transmission of compressed video over OFDM channels," in *Proc. IEEE ICME AVCC*, Barcelona, Spain, Jul. 2011, pp. 1–6.
- [17] N. Golrezai, A. G. Dimakis, and A. F. Molisch, "Device-to-device collaboration through distributed storage," in *Proc. IEEE GLOBECOM*, Anaheim, CA, USA, Dec. 2012, pp. 2397–2402.
- [18] K. Shanmugam and G. Caire, "Wireless downloading delay under proportional fair scheduling with coupled service and requests: An approximated analysis," in *Proc. IEEE ISIT*, Boston, MA, USA, Jul. 2012, pp. 2841–2845.
- [19] N. Golrezai, A. G. Dimakis, and A. F. Molisch, "Wireless device-to-device communications with distributed caching," in *Proc. IEEE ISIT*, Boston, MA, USA, Jul. 2012, pp. 2781–2785.
- [20] N. Golrezai, A. F. Molisch, and A. Dimakis, "Base-station assisted device-to-device communications for high-throughput wireless video networks," in *Proc. IEEE ICC ViOpt*, Ottawa, ON, Canada, Jun. 2012, pp. 7077–7081.
- [21] N. Golrezai, K. Shanmugam, A. Dimakis, A. F. Molisch, and G. Caire, "Wireless video content delivery through coded distributed caching," in *Proc. IEEE ICC*, Ottawa, ON, Canada, Jun. 2012, pp. 2467–2472.
- [22] N. Golrezai, K. Shanmugam, A. Dimakis, A. F. Molisch, and G. Caire, "FemtoCaching: Wireless video content delivery through distributed caching helpers," in *Proc. IEEE INFOCOM*, Orlando, FL, USA, Mar. 2012, pp. 1107–1115.
- [23] N. Golrezai, A. G. Dimakis, A. F. Molisch, and G. Caire, "Wireless video content delivery through distributed caching and peer-to-peer gossiping," in *Proc. Asilomar Conf. Signals, Syst. Comput.*, Pacific Grove, CA, USA, Nov. 2011, pp. 1177–1180.
- [24] M. Ji, G. Caire, and A. F. Molisch, "Wireless device-to-device caching networks: Basic principles and system performance," 2014 [Online]. Available: <http://arxiv.org/abs/1305.5216>
- [25] X. Wu *et al.*, "FlashLinQ: A synchronous distributed scheduler for peer-to-peer Ad Hoc networks," in *Proc. Allerton Conf. Commun., Control, Comput.*, Monticello, IL, USA, Oct. 2010, pp. 514–521.
- [26] X. Wu *et al.*, "FlashLinQ: A synchronous distributed scheduler for peer-to-peer Ad Hoc networks," *IEEE/ACM Trans. Netw.*, vol. 21, no. 4, pp. 1215–1228, Aug. 2013.

- [27] F. Baccelli *et al.*, "On optimizing CSMA for wide area Ad-hoc networks," *Queueing Syst.*, vol. 72, pp. 31–68, 2012.
- [28] X. Wang, M. Chen, T. T. Kwon, L. T. Yang, and V. C. M. Leung, "AMES-cloud: A framework of adaptive mobile video streaming and efficient social video sharing in the clouds," *IEEE Trans. Multimedia*, vol. 15, no. 14, pp. 811–820, Jun. 2013.
- [29] X. Wang, T. T. Kwon, Y. Choi, H. Wang, and J. Liu, "Cloud-assisted adaptive video streaming and social-aware video prefetching for mobile users," *IEEE Wireless Commun.*, vol. 20, no. 3, pp. 72–79, Jun. 2013.
- [30] H. Mansour, V. Krishnamurthy, and P. Nasiopoulos, "Channel aware multiuser scalable video streaming over lossy under-provisioned channels: Modeling and analysis," *IEEE Trans. Multimedia*, vol. 10, no. 7, pp. 1366–1381, Nov. 2008.
- [31] L. de Cicco and S. Mascolo, "An adaptive video streaming control system: Modeling, validation, and performance evaluation," *IEEE/ACM Trans. Netw.*, vol. 22, no. 2, pp. 526–539, Apr. 2014.
- [32] J. Chakareski, J. G. Apostolopoulos, S. Wee, W.-T. Tan, and B. Girod, "Rate-distortion hint tracks for adaptive video streaming," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 15, no. 10, pp. 1257–1269, Oct. 2005.
- [33] S. Tavakoli, J. Gutierrez, and N. Garcia, "Subjective quality study of adaptive streaming of monoscopic and stereoscopic video," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 4, pp. 684–692, Apr. 2014.
- [34] T. C. Thang *et al.*, "Adaptive video streaming over HTTP with dynamic resource estimation," *J. Commun. Netw.*, vol. 15, no. 6, pp. 635–644, Dec. 2013.
- [35] K. Doppler, M. Rinne, C. Wijting, C. B. Ribeiro, and K. Hugl, "Device-to-device communication as an underlay to LTE-advanced networks," *IEEE Commun. Mag.*, vol. 47, no. 12, pp. 42–49, Dec. 2009.
- [36] S. Sanghavi, D. Shah, and A. S. Willsky, "Message passing for maximum weight independent set," in *Proc. NIPS*, Vancouver, BC, Canada, Dec. 2007, pp. 1281–1288.
- [37] S. Sanghavi, D. Shah, and A. S. Willsky, "Message passing for maximum weight independent set," *IEEE Trans. Inf. Theory*, vol. 55, no. 11, pp. 4822–4834, Oct. 2009.
- [38] Y. S. de la Fuente *et al.*, "Efficient HTTP-based streaming using scalable video coding," *Signal Process., Image Commun.*, vol. 27, no. 4, pp. 329–342, Apr. 2012.
- [39] T. Schierl, Y. S. de la Fuente, R. Globisch, C. Hellge, and T. Wiegand, "Priority-based media delivery using SVC with RTP and HTTP streaming," *Multimedia Tools Appl.*, vol. 55, no. 2, pp. 227–246, 2011.
- [40] Y. S. de la Fuente *et al.*, "iDASH: Improved dynamic adaptive streaming over HTTP using scalable video coding," in *Proc. ACM MMSys*, San Jose, CA, USA, Feb. 2011, pp. 257–264.
- [41] W. Wu, R. T. B. Ma, and J. C. S. Lui, "Distributed caching via rewarding: An incentive scheme design in P2P-VoD systems," *IEEE Trans. Parallel Distrib. Syst.*, vol. 25, no. 3, pp. 612–621, Mar. 2014.
- [42] W. Wu, J. C. S. Lui, and R. T. B. Ma, "On incentivizing upload capacity in P2P-VoD systems: Design, analysis and evaluation," *Comput. Netw.*, vol. 57, no. 7, pp. 1674–1688, May 2013.
- [43] W. Wu and J. C. S. Lui, "Exploring the optimal replication strategy in P2P-VoD systems: Characterization and evaluation," *IEEE Trans. Parallel Distrib. Syst.*, vol. 23, no. 8, pp. 1492–1503, Aug. 2012.
- [44] Y. Chen, B. Zhang, C. Chen, and D. M. Chiu, "Performance modeling and evaluation of peer-to-peer live streaming systems under flash crowds," *IEEE/ACM Trans. Netw.*, vol. 22, no. 4, pp. 1106–1120, Aug. 2014.
- [45] Y. Zhou, T. Z. J. Fu, and D. M. Chiu, "On replication algorithms in P2P VoD," *IEEE/ACM Trans. Netw.*, vol. 21, no. 1, pp. 233–243, Feb. 2013.
- [46] Y. Zhou, T. Z. J. Fu, and D. M. Chiu, "Division-of-labor between server and P2P for streaming VoD," in *Proc. IEEE/ACM IWQoS*, Coimbra, Portugal, Jun. 2012, pp. 1–9.
- [47] Y. Liu, Y. Guo, and C. Liang, "A survey on peer-to-peer video streaming systems," *Peer-to-Peer Netw. Appl.*, vol. 1, pp. 18–28, 2008.
- [48] S. Gheorghiu, L. Lima, A. L. Toledo, J. Barros, and M. Medard, "On the performance of network coding in multi-resolution wireless video streaming," in *Proc. IEEE NetCod*, Toronto, ON, Canada, Jun. 9–11, 2010, pp. 1–6.
- [49] M. Feldman, C. Papadimitriou, J. Chuang, and I. Stoica, "Free-riding and whitewashing in peer-to-peer systems," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 5, pp. 1010–1019, May 2006.
- [50] I. Stoica *et al.*, "Chord: A scalable peer-to-peer lookup protocol for internet applications," *IEEE/ACM Trans. Netw.*, vol. 11, no. 1, pp. 17–32, Feb. 2003.
- [51] L. Keller *et al.*, "MicroCast: Cooperative video streaming on smart-phones," in *Proc. ACM MobiSys*, Lake District, U.K., Jun. 25–29, 2012, pp. 57–70.
- [52] S. Ghandeharizadeh, B. Krishnamachari, and S. Song, "Placement of continuous media in wireless peer-to-peer networks," *IEEE Trans. Multimedia*, vol. 6, no. 2, pp. 335–342, Apr. 2004.
- [53] H.-Y. Hsieh and R. Sivakumar, "On using peer-to-peer communication in cellular wireless data networks," *IEEE Trans. Mobile Comput.*, vol. 3, no. 1, pp. 57–72, Jan. 2004.
- [54] D. Bethanabhotla, G. Caire, and M. J. Neely, "Joint transmission scheduling and congestion control for adaptive streaming in wireless device-to-device networks," in *Proc. Asilomar Conf. Signals, Syst., Comput.*, Pacific Grove, CA, USA, Nov. 2012, pp. 1179–1183.
- [55] D. Bethanabhotla, G. Caire, and M. J. Neely, "Adaptive video streaming for wireless networks with multiple users and helpers," *IEEE Trans. Commun.*, vol. 63, no. 1, pp. 268–285, Jan. 2015.
- [56] A. Mishra, M. Shin, and W. Arbaugh, "An empirical analysis of the IEEE 802.11 MAC layer handoff process," *Comput. Commun. Rev.*, vol. 33, no. 2, pp. 93–102, Apr. 2003.
- [57] I. Ramani and S. Savage, "SyncScan: Practical fast handoff for 802.11 infrastructure networks," in *Proc. IEEE INFOCOM*, Miami, FL, USA, Mar. 2005, pp. 675–684.
- [58] L. Lei, Z. Zhong, C. Lin, and X. Shen, "Operator controlled device-to-device communications in LTE-advanced networks," *IEEE Commun. Mag.*, vol. 50, no. 6, pp. 96–104, Jun. 2012.
- [59] S. Winkler and P. Mohandas, "The evolution of video quality measurement: From PSNR to hybrid metrics," *IEEE Trans. Broadcast.*, vol. 54, no. 3, pp. 660–668, Sep. 2008.
- [60] N. Staelens *et al.*, "Constructing a no-reference H.264/AVC bitstream-based video quality metric using genetic programming-based symbolic regression," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 23, no. 8, pp. 1322–1333, Aug. 2013.
- [61] M. J. Neely, *Stochastic Network Optimization with Application to Communication and Queueing Systems*. San Rafael, CA, USA: Morgan & Claypool, 2010.
- [62] M. Z. Win, P. C. Pinto, and L. A. Shepp, "A mathematical theory of network interference and its applications," *Proc. IEEE*, vol. 97, no. 2, pp. 205–230, Feb. 2009.
- [63] A. F. Molisch, *Wireless Communications*, 2nd ed. Piscataway, NJ, USA: IEEE-Wiley, Feb. 2011.
- [64] J. Kim *et al.*, "Demo: Adaptive video streaming for device-to-device mobile platforms," in *Proc. ACM MobiCom*, Miami, FL, USA, Sep. 2013, pp. 127–130.



Joongheon Kim (M'06) received the B.S. and M.S. degrees in computer science and engineering from Korea University, Seoul, Korea, in 2004 and 2006, respectively, and the Ph.D. degree in computer science from the University of Southern California (USC), Los Angeles, CA, USA, in 2014.

Before joining USC, he was a Research Engineer with LG Electronics, Seoul, Korea, from 2006 to 2009. He has been a Systems Engineer with Intel Corporation, Santa Clara, CA, USA, since 2013.

His current research interests are millimeter-wave backhaul/cellular radio platforms and device-to-device streaming platforms.

Dr. Kim is a member of the IEEE Communications Society and IEEE Young Professionals. He received the USC Annenberg Graduate Fellowship Award with his Ph.D. admission from USC in 2009.



Giuseppe Caire (S'92–M'94–SM'03–F'05) was born in Turin, Italy, in 1965. He received the B.Sc. degree from Politecnico di Torino, Turin, Italy, in 1990, the M.Sc. degree from Princeton University, Princeton, NJ, USA, in 1992, and the Ph.D. degree from Politecnico di Torino in 1994, all in electrical engineering.

He has been a Post-Doctoral Research Fellow with the European Space Agency (ESTEC), Noordwijk, The Netherlands, from 1994 to 1995; Assistant Professor in telecommunications with Politecnico di Torino; Associate Professor with the University of Parma, Parma, Italy; and Professor with the Department of Mobile Communications, Eurecom

Institute, Sophia-Antipolis, France. He is currently a Professor of electrical engineering with the Viterbi School of Engineering, University of Southern California, Los Angeles, CA, USA, and an Alexander von Humboldt Professor with the Electrical Engineering and Computer Science Department, Technical University of Berlin, Berlin, Germany. His main research interests are in the fields of communications theory, information theory, and channel and source coding with particular focus on wireless communications.

Prof. Caire has served on the Board of Governors of the IEEE Information Theory Society from 2004 to 2007, and as an officer from 2008 to 2013. He was President of the IEEE Information Theory Society in 2011. He served as Associate Editor for the IEEE TRANSACTIONS ON COMMUNICATIONS from 1998 to 2001 and as Associate Editor for the IEEE TRANSACTIONS ON INFORMATION THEORY from 2001 to 2003. He received the Jack Neubauer Best System Paper Award from the IEEE Vehicular Technology Society in 2003, the IEEE Communications Society and Information Theory Society Joint Paper Award in 2004 and 2011, the Okawa Research Award in 2006, the Alexander von Humboldt Professorship in 2014, and the Vodafone Innovation Prize in 2015.



Andreas F. Molisch (S'89–M'95–SM'00–F'05) received the Dipl. Ing., Ph.D., and habilitation degrees from the Technical University of Vienna, Vienna, Austria, in 1990, 1994, and 1999, respectively.

He subsequently was with AT&T (Bell) Laboratories Research, Middletown, NJ, USA; Lund University, Lund, Sweden; and Mitsubishi Electric Research Labs, Cambridge, MA, USA. He is now a Professor of electrical engineering and Director of the Communication Sciences Institute with the University of Southern California, Los Angeles, CA, USA. He has

authored, coauthored, or edited four books [among them the textbook *Wireless Communications* (Wiley-IEEE Press, 2011)], 16 book chapters, some 170 journal papers, 250 conference papers, as well as more than 80 patents and 70 standards contributions. His current research interests are the measurement and modeling of mobile radio channels, ultrawideband communications and localization, cooperative communications, multiple-input–multiple-output systems, wireless systems for healthcare, and novel cellular architectures.

Dr. Molisch is a Fellow of the AAAS, Fellow of the IET, an IEEE Distinguished Lecturer, and a member of the Austrian Academy of Sciences. He has been an Editor of a number of journals and special issues; General Chair, Technical Program Committee Chair, or Symposium Chair of multiple international conferences; as well as Chairman of various international standardization groups. He has received numerous awards, among them the Donald Fink Prize of the IEEE and the Eric Sumner Award of the IEEE.