

# Cooperative Tx/Rx Caching in Interference Channels: A Storage-Latency Tradeoff Study

Fan Xu, Kangqi Liu and Meixia Tao

Dept. of Electronic Engineering, Shanghai Jiao Tong University, Shanghai, China

Emails: xxiaof@sjtu.edu.cn, k.liu.cn@ieee.org, mxtao@sjtu.edu.cn

**Abstract**—This paper studies the storage-latency tradeoff in the  $3 \times 3$  wireless interference network with caches equipped at all transmitters and receivers. The tradeoff is characterized by the so-called *fractional delivery time* (FDT) at given normalized transmitter and receiver cache sizes. We first propose a generic cooperative transmitter/receiver caching strategy with adjustable file splitting ratios. Based on this caching strategy, we then design the delivery phase carefully to turn the considered interference channel opportunistically into broadcast channel, multicast channel, X channel, or a hybrid form of these channels. After that, we obtain an achievable upper bound of the minimum FDT by solving a linear programming problem of the file splitting ratios. The achievable FDT is a convex and piece-wise linear decreasing function of the cache sizes. Receiver local caching gain, coded multicasting gain, and transmitter cooperation gain (interference alignment and interference neutralization) are leveraged in different cache size regions.

## I. INTRODUCTION

Mobile data traffic has been shifting from connection-centric services, such as voice, e-mails and web browsing, to emerging content-centric services, such as video streaming, push media, application download/updates, and mobile TV. The contents in these services are typically produced well ahead of transmission and can be requested by multiple users, although at possibly different times. This allows us to cache the contents at the edge of wireless networks, e.g., base stations and user devices, and hence to reduce user access latency and alleviate wireless traffic. A fundamental question in wireless cache networks is what and how much gain can be leveraged through caching.

Caching in a shared link with one server and multiple cache-enabled users is first studied by Maddah-Ali and Niesen in [1]. It is shown that caching at user ends brings not only local caching gain but also global caching gain. The latter is achieved by a carefully designed cache placement and coded delivery strategy, which can create multicast chances for content delivery even if users demand different files. The idea is then extended to the decentralized coded caching in a large network in [2]. In [3], the authors considered the wireless broadcast channel with imperfect channel state information at the transmitter (CSIT) and showed that the gain of coded multicasting scheme can offset the loss due to the imperfect CSIT.

This work is supported by the NSF of China under grants 61571299, 61322102, and 61329101.

The authors in [4] studied the transmitter cache strategy in the cache-aided interference channel. It is shown that splitting contents into different parts and caching each part in different transmitters can turn the interference channel into broadcast channel, X channel, or hybrid channel and hence increase the system throughput via interference management. The authors in [5] presented a lower bound of delivery latency in a general interference network with transmitter cache and showed that the scheme in [4] is optimal in certain region of cache size.

The above literature reveals that caching at the receiver side can bring coded multicasting gain, and that caching at the transmitter side can induce transmitter cooperation for interference management. It is thus of great interests to investigate the impact of caching at both transmitter and receiver sides.

In this paper, we aim to study the fundamental limits of caching in the  $3 \times 3$  interference network with caches equipped at all transmitters and receivers as shown in Fig. 1. We adopt the storage-latency tradeoff originally proposed in [5] to characterize the fundamental limits. In specific, we measure the performance by the *fractional delivery time* (FDT) as a function of the normalized receiver and transmitter cache sizes. To analyze the minimum FDT, we propose a generic file splitting and caching strategy with adjustable file splitting ratios. Based on this strategy, we then design the delivery phase carefully so that the network topology can be opportunistically changed to broadcast channel, multicast channel, X channel, or a hybrid form of these channels. We then obtain an achievable upper bound of the minimum FDT by optimizing the file splitting ratios. The obtained FDT is a convex and piece-wise linear decreasing function of the transmitter and receiver cache sizes. Our result shows that coded multicasting gain should be exploited as much as we can when the cache sizes are very limited. It also shows that transmitter cooperation gain can only be exploited when the transmitter cache size exceeds a certain threshold dependent on the receiver cache size. Note that an independent work on the similar problem is studied in [6]. We shall discuss the differences with [6] later.

Notations:  $(\cdot)^T$  denotes the transpose.  $[K]$  denotes set  $\{1, 2, \dots, K\}$ .  $\lfloor x \rfloor$  denotes the largest integer no greater than  $x$ .  $(x_j)_{j=1}^K$  denotes vector  $(x_1, x_2, \dots, x_K)^T$ .  $\mathcal{CN}(m, \sigma^2)$  denotes the complex Gaussian distribution with mean of  $m$  and variance of  $\sigma$ .

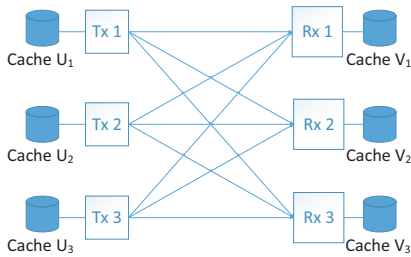


Fig. 1:  $3 \times 3$  Interference channel with cache at Tx/Rx sides.

## II. SYSTEM MODEL AND DEFINITIONS

Consider the  $3 \times 3$  cache-aided interference channel shown in Fig. 1. Each node is assumed to have single antenna. The communication link between each transmitter and each receiver experiences channel fading and is corrupted with additive white Gaussian noise. The communication at each time slot  $t$  over this network is modeled by

$$Y_j(t) = \sum_{p=1}^3 h_{jp}(t)X_p(t) + Z_j(t), j = 1, 2, 3,$$

where  $Y_j(t) \in \mathbb{C}$  denotes the received signal at receiver  $j$ ,  $X_p(t) \in \mathbb{C}$  denotes the transmitted signal at transmitter  $p$ ,  $h_{jp}(t) \in \mathbb{C}$  denotes the channel coefficient from transmitter  $p$  to receiver  $j$ , and  $Z_j(t)$  denotes the noise at receiver  $j$  distributed as  $\mathcal{CN}(0, 1)$ .

Consider a database consisting of  $L$  files ( $L \gg 3$ ), denoted by  $\{W_1, W_2, \dots, W_L\}$ . Each file is chosen independently and uniformly from  $[2^F] = \{1, 2, \dots, 2^F\}$  randomly, where  $F$  is the file size in bits. Each transmitter has a local cache able to store  $M_T F$  bits and each receiver has a local cache able to store  $M_R F$  bits. The *normalized cache sizes* at each transmitter and receiver are defined, respectively, as

$$\mu_T \triangleq \frac{M_T}{L}, \quad \mu_R \triangleq \frac{M_R}{L}.$$

The network operates in two phases, *cache placement phase* and *content delivery phase*. During the cache placement phase, each transmitter  $p$  designs a caching function

$$\phi_p : [2^F]^L \rightarrow [2^{FM_T}],$$

mapping the  $L$  files in the database to its local cached content  $U_p \triangleq \phi_p(W_1, W_2, \dots, W_L)$ . Each receiver  $j$  also designs a caching function

$$\psi_j : [2^F]^L \rightarrow [2^{FM_R}],$$

mapping the  $L$  files to its local cached content  $V_j \triangleq \psi_j(W_1, W_2, \dots, W_L)$ . The caching functions  $\{\phi_p, \psi_j\}$  are assumed to be known globally at all nodes. In the delivery phase, each receiver  $j$  requests a file  $W_{d_j}$  from the database. We denote  $\mathbf{d} \triangleq (d_j)_{j=1}^3 \in [L]^3$  as the demand vector. Each transmitter  $p$  has an encoding function

$$\Lambda_p : [2^{FM_T}] \times [L]^3 \times \mathbb{C}^{3 \times 3} \rightarrow \mathbb{C}^T.$$

Transmitter  $p$  uses  $\Lambda_p$  to map its cached content  $U_p$ , receiver demands  $\mathbf{d}$  and channel realization  $\mathbf{H}$  to the signal

$(X_p[t])_{t=1}^T \triangleq \Lambda_p(U_p, \mathbf{d}, \mathbf{H})$ , where  $T$  is the block length of the code. Note that  $T$  may depend on the receiver demand  $\mathbf{d}$  and channel realization  $\mathbf{H}$ , and thus can also be denoted as  $T^{\mathbf{d}, \mathbf{H}}$  (with a slight abuse of notation, we will use  $T$  again to denote the average worst-case delivery time in Definition 1). Each codeword  $(X_p[t])_{t=1}^T$  has an average transmit power constraint  $P$ . Each receiver  $j$  has a decoding function

$$\Gamma_j : [2^{FM_R}] \times \mathbb{C}^T \times \mathbb{C}^{3 \times 3} \times [L]^3 \rightarrow [2^F].$$

Upon receiving  $(Y_j[t])_{t=1}^T$ , each receiver  $j$  uses  $\Gamma_j$  to decode  $\hat{W}_j \triangleq \Gamma_j(V_j, (Y_j[t])_{t=1}^T, \mathbf{H}, \mathbf{d})$  of its desired file  $W_{d_j}$  using its cached content  $V_j$  as side information. The worst-case error probability is

$$P_\epsilon = \max_{\mathbf{d} \in [L]^3} \max_{j \in [3]} \mathbb{P}(\hat{W}_j \neq W_{d_j}).$$

The given caching and coding scheme  $\{\phi_p, \psi_j, \Lambda_p, \Gamma_j\}$  is said to be feasible if  $P_\epsilon \rightarrow 0$  when  $F \rightarrow \infty$ .

In this work, we adopt the following latency-oriented performance metrics originally proposed in [5].

**Definition 1:** The delivery time (DT) for a given feasible caching and coding scheme is defined as

$$T \triangleq \lim_{P \rightarrow \infty} \lim_{F \rightarrow \infty} \max_{\mathbf{d}} \mathbb{E}_{\mathbf{H}}(T^{\mathbf{d}, \mathbf{H}}). \quad (1)$$

**Definition 2:** The *fractional delivery time* (FDT) for a given feasible caching and coding scheme is defined as

$$\tau(\mu_R, \mu_T) \triangleq \lim_{P \rightarrow \infty} \lim_{F \rightarrow \infty} \sup_{\mathbf{d}} \frac{\max_{\mathbf{d}} \mathbb{E}_{\mathbf{H}}(T^{\mathbf{d}, \mathbf{H}})}{N_R F \cdot 1/\log P},$$

where  $N_R = 3$  is the number of content requesters. Moreover, the minimum FDT at given normalized cache sizes  $\mu_T$  and  $\mu_R$  is defined as

$$\tau^*(\mu_R, \mu_T) = \inf\{\tau(\mu_R, \mu_T) : \tau(\mu_R, \mu_T) \text{ is achievable}\}.$$

The above performance metrics are defined in the asymptotic sense when  $P \rightarrow \infty$  and  $F \rightarrow \infty$ . It is clear that the FDT and DT are related by  $\tau = \frac{T \log P}{3F}$ . The FDT  $\tau(\mu_R, \mu_T)$  can be regarded as the relative time with respect to delivering the total  $3F$  requested bits in an interference-free baseline system with transmission rate  $\log P$  in the high SNR region.

**Remark 1:** Our definition of FDT  $\tau$  is slightly different from the *normalized delivery time* (NDT)  $\delta$  in [5] in that our FDT is further normalized by the number of receivers. That is,  $\tau = \delta/3$ . With such normalization, the FDT is defined for the total  $3F$  bits requested in the network rather than the  $F$  bits requested by a single receiver as in [5]. As a result, the range of FDT is  $0 \leq \tau \leq 1$ , which is truly normalized.

**Remark 2:** Compared to the ‘‘load’’  $R$  defined for the shared link in [1], the FDT can be expressed as  $\tau = \frac{R}{3 \cdot \text{DoF}}$ , where DoF is the sum DoF of the considered channel. Comparing to the standard DoF adopted for interference channel with transmitter cache only in [4], we have  $\tau(\mu_R = 0, \mu_T) = \frac{1}{\text{DoF}}$ . As a result, the FDT evaluates the delivery time of the actual *load* at a transmission rate specified by *DoF* of the given channel, and hence is particularly suitable to characterize the performance

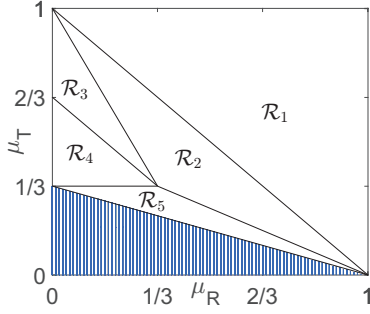


Fig. 2: Feasible domain of FDT (divided into 5 regions).

of the wireless network with both transmitter and receiver caches.

**Remark 3** (Feasible domain of FDT): The FDT introduced above is able to measure the fundamental tradeoff between the cache storage and content delivery latency. However, not all normalized cache sizes are feasible. Given fixed  $L$  and  $M_T$ , all the transmitters together can store at most  $3M_T F$  bits of files, which leaves  $LF - 3M_T F$  bits of files to be stored in all receivers. Thus we must have  $M_R F \geq LF - 3M_T F$ . This gives the feasible region for the normalized cache sizes as:

$$\begin{cases} 0 \leq \mu_R, \mu_T \leq 1 \\ \mu_R + 3\mu_T \geq 1 \end{cases}. \quad (2)$$

Throughout this paper, we study the FDT only in the feasible domain (2).

### III. MAIN RESULTS

In this section, we present an achievable upper bound of the minimum FDT  $\tau^*(\mu_R, \mu_T)$ . The proof will be given in the next two sections.

**Theorem 1:** For the  $3 \times 3$  cache-aided interference channel, the minimum FDT is upper bounded by

$$\tau^*(\mu_R, \mu_T) \leq \begin{cases} \frac{1}{3} - \frac{\mu_R}{3}, & (\mu_R, \mu_T) \in \mathcal{R}_1 \\ \frac{4}{9} - \frac{4\mu_R}{9} - \frac{\mu_T}{9}, & (\mu_R, \mu_T) \in \mathcal{R}_2 \\ \frac{1}{2} - \frac{5}{9}\mu_R - \frac{1}{6}\mu_T, & (\mu_R, \mu_T) \in \mathcal{R}_3 \\ \frac{13}{18} - \frac{8}{9}\mu_R - \frac{1}{2}\mu_T, & (\mu_R, \mu_T) \in \mathcal{R}_4 \\ \frac{8}{9} - \frac{8}{9}\mu_R - \mu_T, & (\mu_R, \mu_T) \in \mathcal{R}_5 \end{cases},$$

where  $\{\mathcal{R}_i\}_{i=1}^5$  are given below and sketched in Fig. 2.

$$\begin{cases} \mathcal{R}_1 = \{(\mu_R, \mu_T) : \mu_R + \mu_T \geq 1, \mu_R \leq 1, \mu_T \leq 1\} \\ \mathcal{R}_2 = \{(\mu_R, \mu_T) : \mu_R + \mu_T < 1, 2\mu_R + \mu_T \geq 1, \\ \mu_R + 2\mu_T > 1\} \\ \mathcal{R}_3 = \{(\mu_R, \mu_T) : \mu_R + \mu_T \geq \frac{2}{3}, 2\mu_R + \mu_T < 1, \\ \mu_R \geq 0\} \\ \mathcal{R}_4 = \{(\mu_R, \mu_T) : \mu_R + \mu_T < \frac{2}{3}, \mu_R \geq 0, \mu_T > \frac{1}{3}\} \\ \mathcal{R}_5 = \{(\mu_R, \mu_T) : \mu_T \leq \frac{1}{3}, \mu_R + 2\mu_T \leq 1, \\ \mu_R + 3\mu_T \geq 1\} \end{cases}.$$

The above theorem shows that the achievable FDT is a convex and piecewise linear decreasing function of  $\mu_R$  and  $\mu_T$ . It captures an achievable tradeoff between the cache storage and the delivery latency. In the special case when  $\mu_R = 0$  (transmitters cache only), the results reduce to

$$\tau^*(0, \mu_T) \leq \begin{cases} 13/18 - \mu_T/2, & 1/3 \leq \mu_T \leq 2/3 \\ 1/2 - \mu_T/6, & 2/3 < \mu_T \leq 1 \end{cases},$$

which is the same as the upper bound of  $1/\text{DoF}$  in [4].

When  $\mu_T = 1$ , each transmitter can cache all the files and hence the network can be viewed as a virtual broadcast channel as in [1] except that the transmitter has 3 distributed antennas. Thus, we can achieve FDT  $\tau = \frac{1}{3} - \frac{\mu_R}{3}$  here. Comparing to the result in [1], i.e.,  $\tau = \frac{1-\mu_R}{1+3\mu_R}$  at  $\mu_R = \{0, \frac{1}{3}, \frac{2}{3}, 1\}$ , we can see that our FDT is better when  $0 \leq \mu_R < \frac{2}{3}$  and they are same when  $\frac{2}{3} \leq \mu_R \leq 1$ . The performance improvement is due to transmitter cooperation gain.

## IV. ACHIEVABLE CACHING SCHEME

### A. File Splitting and Placement

Given fixed  $\mu_R$  and  $\mu_T$ , the content placement can be established as follows.

In this work, we treat all the files equally without taking file popularity into account. Thus, each file will be split and cached in the same manner. Without loss of generality, we focus on the splitting and caching of file  $W_i$  for any  $1 \leq i \leq L$ . Since each bit of the file is either cached or not cached in every node, there are  $2^6 = 64$  possible cache states for each bit. However, note that every bit of the file must be cached in at least one node. In addition, every bit that is not cached simultaneously in all receivers must be cached in at least one transmitter. This is because we do not allow receiver cooperation and all the messages must be sent from the transmitters. As such, the total number of feasible cache states for each bit is given by  $64 - 1 - \binom{3}{1} - \binom{3}{2} = 57$ . Now with possibly different lengths, we can partition each  $W_i$  into 57 subfiles exclusively.

Define receiver subset  $\Phi \subseteq [3]$  and transmitter subset  $\Psi \subseteq [3]$ . Then, denote  $W_{i\Phi t\Psi}$  as the subfile of  $W_i$  cached in receivers in  $\Phi$  and transmitters in  $\Psi$ . For example,  $W_{i r_1 t_1}$  is the subfile cached in receiver 1 and transmitter 1,  $W_{i r_0 t_{123}}$  is the subfile cached in none of the three receivers but in three transmitters. Similarly, we denote  $W_{i r\Phi}$  as the collection of the subfiles in file  $W_i$  that are cached in receivers in  $\Phi$ , i.e.  $W_{i r\Phi} = \bigcup_{\Psi} W_{i r\Phi t\Psi}$ . We further assume that the subfiles that are cached in the same number of transmitters and the same number of receivers have the same length. Due to the symmetry of all the nodes as well as the independency of all files, this assumption is valid. Thus, we assume the size of  $W_{i r\Phi t\Psi}$  is  $a_{|\Phi||\Psi|} F$ , where  $|\Psi|$  and  $|\Phi|$  are the cardinalities of  $\Psi$  and  $\Phi$ , respectively, and  $a_{|\Phi||\Psi|}$  is the file splitting ratio to be optimized later. For example, the size of  $W_{i r_1 t_1}$  is  $a_{11} F$  and the size of  $W_{i r_0 t_{123}}$  is  $a_{03} F$ . Here, the file splitting ratios  $\{a_{|\Phi||\Psi|}\}$  should satisfy the following constraints:

$$a_{30} + 3a_{31} + 3a_{32} + a_{33} + 9a_{21} + 9a_{22} + 3a_{23} + 9a_{11} + 9a_{12} + 3a_{13} + 3a_{01} + 3a_{02} + a_{03} = 1, \quad (3)$$

$$a_{30} + 3a_{31} + 3a_{32} + a_{33} + 6a_{21} + 6a_{22} + 2a_{23} + 3a_{11} + 3a_{12} + a_{13} \leq \mu_R, \quad (4)$$

$$a_{31} + 2a_{32} + a_{33} + 3a_{21} + 6a_{22} + 3a_{23} + 3a_{11} + 6a_{12} + 3a_{13} + a_{01} + 2a_{02} + a_{03} \leq \mu_T. \quad (5)$$

Constraint (3) comes from the constraint of file size. The multiplier of each splitting ratio  $a_{|\Phi||\Psi|}$  in (3) indicates the

number of subfiles that have the same length of  $a_{|\Phi||\Psi|}f$ . For instance, the number of subfiles with length  $a_{21}F$  is nine since there are  $\binom{3}{2}\binom{3}{1} = 9$  cache states to cache the subfile in two out of the three receivers and one out of the three transmitters. Constraints (4) and (5) come from the receiver and transmitter cache size limit, respectively. Similar arguments used in (3) can be applied here to determine the multipliers.

### B. File Delivery

Without loss of generality, we assume that receivers 1, 2, 3 desire files  $W_1, W_2, W_3$ , respectively. Specifically, receiver  $j$  ( $j \in [3]$ ) desires subfiles  $W_{jr_0}, W_{jr_{kl}}$ , and  $W_{jr_k}$  that are not cached in its local cache, where  $k, l \neq j$ . We divide these subfiles into three groups and present the delivery scheme for each group individually.

1) *Delivery of Subfiles Cached in Two Receivers:* Consider the delivery of subfiles  $\{W_{jr_{kl}}\}_{k,l \neq j}$  needed by receiver  $j$  ( $j \in [3]$ ). Since the subfiles desired by each receiver are cached in the other two receivers, coded multicasting opportunities can be exploited. In specific, consider subfiles  $W_{1r_{23t_\Psi}}, W_{2r_{13t_\Psi}}$ , and  $W_{3r_{12t_\Psi}}$  for any transmitter subset  $\Psi$ . Transmitters in each subset  $\Psi$  can generate a new message  $W_{123t_\Psi}^\oplus \triangleq W_{1r_{23t_\Psi}} \oplus W_{2r_{13t_\Psi}} \oplus W_{3r_{12t_\Psi}}$  needed by all three receivers, where  $\oplus$  denotes the bit-wise XOR. To illustrate the delivery scheme, we take the set of messages with  $|\Psi| = 2$  for example. The message flow pattern is shown in Fig. 3. We adopt time division multiple access (TDMA) technique so that all the  $\binom{3}{2}$  possible transmitter cooperation sets take turns to transmit. In specific, we divide the transmission time into 3 time slots. In each time slot, we select one transmitter subset  $\Psi$  (e.g.  $\Psi = \{1, 2\}$ ) and let transmitters in this subset to cooperatively transmit  $W_{123t_\Psi}^\oplus$  to all three receivers. The network topology in each slot now becomes a broadcast channel with common information only, which we refer to as *multicast channel*. The maximum sum DoF of the multicast channel is 1, no matter the transmitters cooperate or not. The converse can be proved easily using cut-set bound on each receiver. Thus, the delivery time  $T = \frac{3a_{22}F}{\log P}$  is achieved. In the general case with  $|\Psi| = i$ , we also use the TDMA technique so that all the  $\binom{3}{i}$  possible transmitter cooperation sets take turns to transmit. The delivery time is  $T = \frac{\binom{3}{i}a_{2i}F}{\log P}$ . Thus, the total FDT of subfiles  $\{W_{jr_{kl}}\}_{k,l \neq j}$  is  $\tau = \frac{1}{3} \sum_{i=1}^3 \binom{3}{i} a_{2i}$ .

2) *Delivery of Subfiles Cached in One Receiver:* Consider the delivery of subfiles  $\{W_{jr_k}\}_{k \neq j}$  needed by receiver  $j$  ( $j \in [3]$ ). Since each subfile requested by one receiver is already cached in another receiver, coded multicasting gain can be exploited again. In specific, transmitters in each subset  $\Psi$  can generate a new message  $W_{jkt_\Psi}^\oplus \triangleq W_{jr_{kt_\Psi}} \oplus W_{kr_{jt_\Psi}}$  needed by receivers  $j$  and  $k$ , where  $j, k \in [3], j < k$ .

We first consider the delivery of messages  $\{W_{jkt_\Psi}^\oplus\}_{j < k}$  with  $|\Psi| = 1$ . The message flow pattern is shown in Fig. 4(a), and the network topology can be seen as the hybrid X-multicast channel. Lemma 1 below presents the sum DoF of this channel. The proof is based on interference alignment and given in [7, Appendix B].

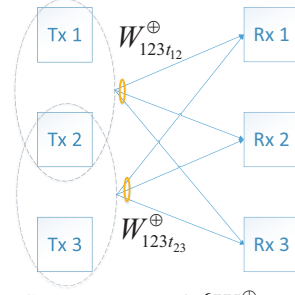


Fig. 3: Message flow pattern of  $\{W_{123t_\Psi}^\oplus\}_{|\Psi|=2}$ . Only  $\Psi = \{1, 2\}$  and  $\{2, 3\}$  are shown. Subfiles are listed beside the channel which carries them. Dashed circle denotes that the transmitters inside it cooperate with each other in the delivery phase. Solid circle denotes that the channels inside it carry the same subfile.

**Lemma 1:** The achievable sum DoF of the  $3 \times 3$  hybrid X-multicast channel is  $\frac{9}{7}$ .

Using Lemma 1 and given that the total amount of bits to deliver is  $9a_{11}F$ , we obtain  $\tau = \frac{7a_{11}}{3}$ .

Next, we consider the delivery of messages  $\{W_{jkt_\Psi}^\oplus\}_{j < k}$  with  $|\Psi| \geq 2$ . The message flow patterns for  $|\Psi| = 2$  and  $|\Psi| = 3$  are shown in Fig. 4(b) and 4(c), where the network topologies can be seen as the partially and fully cooperative hybrid X-multicast channel, respectively. Lemma 2 below presents the sum DoF of this channel. Its proof is based on interference neutralization and given in [7, Appendix C].

**Lemma 2:** The achievable sum DoF of the  $3 \times 3$  partially or fully cooperative hybrid X-multicast channel is  $\frac{3}{2}$ .

As such, the total FDT of subfiles  $\{W_{jr_{kt_\Psi}}\}_{j,k \in [3], j < k}$  for  $\Psi$ 's with  $|\Psi| = 2, 3$  is  $\tau = \frac{6a_{12} + 2a_{13}}{3}$ .

3) *Delivery of Subfiles Cached in None Of Receivers:* Consider the delivery of subfiles  $\{W_{jr_0}\}$  needed by receiver  $j$  ( $j \in [3]$ ). Each  $W_{jr_0}$  further consists of subfiles  $W_{jr_0t_\Psi}$  for all transmitter subsets  $\Psi$ 's with  $|\Psi| = 1, 2, 3$ . The message flow patterns of  $\{W_{jr_0t_\Psi}\}_{|\Psi|=3}$ ,  $\{W_{jr_0t_\Psi}\}_{|\Psi|=2}$  and  $\{W_{jr_0t_\Psi}\}_{|\Psi|=1}$  correspond to the patterns in [4] when  $\mu_T = 1, \frac{2}{3}, \frac{1}{3}$ , respectively. In [4], the message flow patterns of  $\{W_{jr_0t_\Psi}\}_{|\Psi|=3}$ ,  $\{W_{jr_0t_\Psi}\}_{|\Psi|=2}$ , and  $\{W_{jr_0t_\Psi}\}_{|\Psi|=1}$  form a MISO broadcast channel, a partially cooperative X channel, and an X channel, respectively. Thus, the delivery time of subfiles  $\{W_{jr_0t_\Psi}\}_{j \in [3]}$  for all  $\Psi$ 's is  $T = \frac{3a_{03}F}{3 \log P} + \frac{9a_{02}F}{18 \log P/7} + \frac{9a_{01}F}{9 \log P/5}$  and its corresponding FDT is  $\tau = \frac{a_{03}}{3} + \frac{7a_{02}}{6} + \frac{5a_{01}}{3}$ .

## V. CACHING OPTIMIZATION

Combining all the FDTs obtained in Section IV-B, we obtain the total FDT in the delivery phase as

$$\tau = \frac{1}{3}(5a_{01} + \frac{7}{2}a_{02} + a_{03} + 3a_{21} + 3a_{22} + a_{23} + 7a_{11} + 6a_{12} + 2a_{13}). \quad (6)$$

Our goal is to minimize the FDT subject to the file slitting ratio constraints (3)(4)(5). This is formulated as:

$$\begin{aligned} \min \quad & \tau(\mu_R, \mu_T) \\ \text{s.t.} \quad & (3)(4)(5), \end{aligned} \quad (7)$$



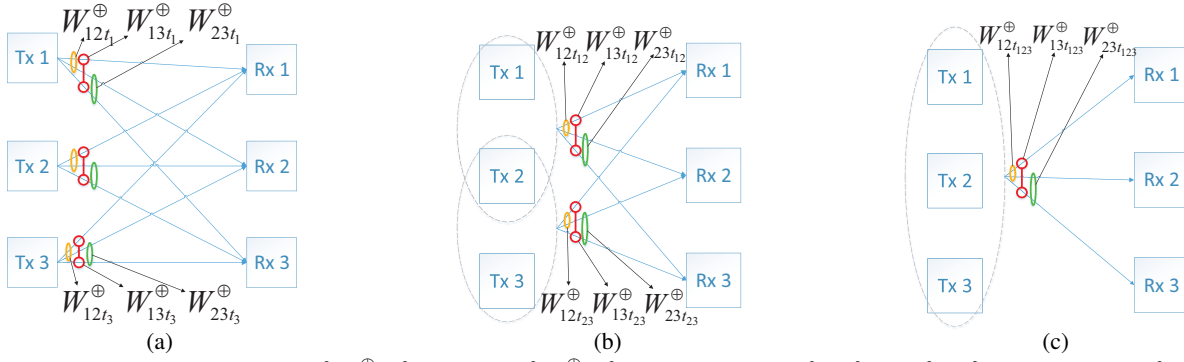


Fig. 4: Message flow pattern of (a)  $\{W_{jkt_\Psi}^\oplus\}_{|\Psi|=1}$ , (b)  $\{W_{jkt_\Psi}^\oplus\}_{|\Psi|=2}$ , only  $\Psi = \{1, 2\}$  and  $\{2, 3\}$  are shown, (c)  $\{W_{jkt_\Psi}^\oplus\}_{|\Psi|=3}$ .

which is a standard linear programming problem. Using linear equation substitution and other manipulations, we can obtain the optimal solutions in closed form as follows. Here, all the regions are defined in Theorem 1.

**Region  $\mathcal{R}_1$ :** The optimal FDT is  $\tau^* = \frac{1}{3} - \frac{\mu_R}{3}$ . The optimal splitting ratios are not unique but must satisfy that  $a_{11}^* = a_{01}^* = a_{02}^* = 0$  and that the equality in (4) holds. One feasible solution is  $a_{30}^* = \mu_R$ ,  $a_{03}^* = 1 - \mu_R$  and other ratios are 0.

**Region  $\mathcal{R}_2$ :** The optimal FDT is  $\tau^* = \frac{4}{9} - \frac{4\mu_R}{9} - \frac{\mu_T}{9}$ . The optimal splitting ratios are not unique but must satisfy

$$\begin{cases} a_{01}^* = a_{02}^* = a_{31}^* = a_{32}^* = a_{33}^* = a_{22}^* = a_{23}^* = a_{13}^* = 0 \\ a_{11}^* = \frac{1}{3} - \frac{\mu_R}{3} - \frac{\mu_T}{3} \\ a_{30}^* + 6a_{21}^* + 3a_{12}^* = 2\mu_R + \mu_T - 1 \\ 3a_{21}^* + 6a_{12}^* + a_{03}^* = \mu_R + 2\mu_T - 1 \end{cases}.$$

One feasible solution is  $a_{11}^* = \frac{1}{3} - \frac{\mu_R}{3} - \frac{\mu_T}{3}$ ,  $a_{30}^* = 2\mu_R + \mu_T - 1$ ,  $a_{03}^* = \mu_R + 2\mu_T - 1$  and other ratios are 0.

**Region  $\mathcal{R}_3$ :** The optimal FDT is  $\tau^* = \frac{1}{2} - \frac{5}{9}\mu_R - \frac{1}{6}\mu_T$ . The optimal splitting ratios are unique and given by  $a_{11}^* = \frac{\mu_R}{3}$ ,  $a_{02}^* = 1 - 2\mu_R - \mu_T$ ,  $a_{03}^* = 3\mu_R + 3\mu_T - 2$  and other ratios being 0.

**Region  $\mathcal{R}_4$ :** The optimal FDT is  $\tau^* = \frac{13}{18} - \frac{8}{9}\mu_R - \frac{1}{2}\mu_T$ . The optimal splitting ratios are unique and given by  $a_{11}^* = \frac{\mu_R}{3}$ ,  $a_{01}^* = \frac{2}{3} - \mu_R - \mu_T$ ,  $a_{02}^* = \mu_T - \frac{1}{3}$  and other ratios being 0.

**Region  $\mathcal{R}_5$ :** The optimal FDT is  $\tau^* = \frac{8}{9} - \frac{8}{9}\mu_R - \mu_T$ . The optimal splitting ratios are unique and given by  $a_{11}^* = \frac{\mu_R}{3} + \mu_T - \frac{1}{3}$ ,  $a_{01}^* = 1 - \mu_R - 2\mu_T$ ,  $a_{30}^* = 1 - 3\mu_T$  and other ratios being 0.

Summarizing all the results above, we finish proof of Theorem 1.

**Remark 4:** In  $\mathcal{R}_1$  and  $\mathcal{R}_2$ , the multiple choices of file splitting ratios from caching optimization offer freedom to choose appropriate caching and delivery scheme in practical systems according to different limitations, such as file splitting constraints.

**Remark 5:** In the proposed caching strategy, the local caching gain, transmitter cooperation gain, and coded multicasting gain are exploited opportunistically in different cache size regions. These gains are reflected by the file splitting ratios of the corresponding cache states. In  $\mathcal{R}_1$ , local caching gain and cooperation gain are exploited, because the feasible solution is  $a_{30}^* = \mu_R$ ,  $a_{03}^* = 1 - \mu_R$ . We do not need to use up the total cache storage at transmitters. In  $\mathcal{R}_2$ ,  $\mathcal{R}_3$ ,

and  $\mathcal{R}_4$ , all the three gains are exploited since their optimal solutions all satisfy  $a_{11}^* > 0$  and  $a_{02}^* + a_{03}^* > 0$ . In  $\mathcal{R}_3$  and  $\mathcal{R}_4$ , there do not exist two receivers which cache the same content. Instead, each receiver uses up its total cache size to cache the content already cached in only one transmitter, i.e.  $W_{ir_j t_p}$ , to fully exploit coded multicasting gain. In  $\mathcal{R}_5$ , only local caching gain and coded multicasting gain are exploited and no transmitter cooperation can be exploited since the optimal solution satisfies  $a_{mn}^* = 0$ ,  $\forall n \geq 2$ . This is due to that the transmitter cache size is approaching its lower limit  $\mu_T \geq \frac{1}{3}(1 - \mu_R)$  in (2).

**Remark 6:** Although the similar caching problem is considered in [6], their performance metric, caching scheme, and conclusion are significantly different from ours. First, we adopt the FDT as the performance metric, while [6] used the standard DoF. From Remark 2, FDT reflects not only the load reduction due to receiver cache but also the DoF enhancement due to transmitter cache, while the DoF alone cannot reflect the former one. Also, at each given  $(\mu_R, \mu_T)$ , the file splitting ratios in [6] are pre-determined, while our file splitting ratios are obtained by solving a linear programming problem and thus are probably optimal under the given caching strategy. Another difference is that the transmission scheme in [6] is restricted to one-shot linear processing, while we allow interference alignment which may require infinite symbol extension.

## REFERENCES

- [1] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. on Information Theory*, vol. 60, no. 5, pp. 2856–2867, May 2014.
- [2] M. Maddah-Ali and U. Niesen, "Decentralized coded caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Trans. on Networking*, Aug 2015.
- [3] J. Zhang and P. Elia, "Fundamental limits of cache-aided wireless BC: interplay of coded-caching and CSIT feedback," 2015. [Online]. Available: <http://arxiv.org/abs/1511.03961>
- [4] M. Maddah-Ali and U. Niesen, "Cache-aided interference channels," in *IEEE International Symposium on Information Theory (ISIT)*, June 2015.
- [5] A. Sengupta, R. Tandon, and O. Simeone, "Cache aided wireless networks: Tradeoffs between storage and latency," 2015. [Online]. Available: <http://arxiv.org/abs/1512.07856>
- [6] N. Naderializadeh, M. A. Maddah-Ali, and A. S. Avestimehr, "Fundamental limits of cache-aided interference management," 2016. [Online]. Available: <http://arxiv.org/abs/1602.04207>
- [7] F. Xu, M. Tao, and K. Liu, "Fundamental tradeoff between storage and latency in cache-aided wireless interference networks," 2016. [Online]. Available: <http://arxiv.org/abs/1605.00203>