

Fundamental Limits of Cache-Aided Interference Management

Navid Naderializadeh*, Mohammad Ali Maddah-Ali[†], and A. Salman Avestimehr*

*Department of Electrical Engineering, University of Southern California, Los Angeles, CA, USA

[†]Nokia Bell Labs, Holmdel, NJ, USA

E-mails: naderial@usc.edu, mohammad.maddah-ali@nokia.com, avestimehr@ee.usc.edu

Abstract

We consider a system, comprising a library of N files (e.g., movies) and a wireless network with K_T transmitters, each equipped with a local cache of size of M_T files, and K_R receivers, each equipped with a local cache of size of M_R files. Each receiver will ask for one of the N files in the library, which needs to be delivered. The objective is to design the cache placement (without prior knowledge of receivers' future requests) and the communication scheme to maximize the throughput of the delivery. In this setting, we show that the sum degrees-of-freedom (sum-DoF) of $\min \left\{ \frac{K_T M_T + K_R M_R}{N}, K_R \right\}$ is achievable, and this is within a factor of 2 of the optimum, under one-shot linear schemes. This result shows that (i) the one-shot sum-DoF scales *linearly* with the *aggregate cache size in the network* (i.e., the cumulative memory available at *all nodes*), (ii) the transmitters' caches and receivers' caches contribute equally in the one-shot sum-DoF, and (iii) caching can offer a throughput gain that scales linearly with the size of the network.

To prove the result, we propose an achievable scheme that exploits the redundancy of the content at transmitters' caches to cooperatively zero-force some outgoing interference, and availability of the unintended content at the receivers' caches to cancel (subtract) some of the incoming interference. We develop a particular pattern for cache placement that maximizes the overall gains of cache-aided transmit and receive interference cancellations. For the converse, we present an integer optimization problem which minimizes the number of communication blocks needed to deliver any set of requested files to the receivers. We then provide a lower bound on the value of this optimization problem, hence leading to an upper bound on the linear one-shot sum-DoF of the network, which is within a factor of 2 of the achievable sum-DoF.

I. INTRODUCTION

Over the last decade, video delivery has emerged as the main driving factor of the wireless traffic. In this context, there is often a large library of pre-recorded content (e.g. movies), out of which, users may request to receive a specific file. One way to reduce the burden of this traffic is to employ memories distributed across the networks and closer to the end users to prefetch some of the popular content. This can help system to deliver the content with higher throughput and less delay.

As a result, there have been significant interests in both academia and industry in characterizing the impact of caching on the performance of communication networks (see, e.g. [1–13]). In particular, in a network with only one transmitter broadcasting to several receivers, it was shown in [2] that local delivery attains only a small fraction of the gain that caching can offer, and by designing a particular pattern in cache placement at the users and exploiting coding in delivery, a significantly larger *global* throughput gain can be achieved, which is a function of the entire cache throughout the network. This also demonstrates that the gain of caching scales with the size of the network. As a follow-up, this work has been extended to the case of multiple transmitters in [3], where it was shown that the gain of caching can be improved if several transmitters have access to the entire library of files. Caching at the transmitters was also considered in [4, 5] and used to induce collaboration between transmitters in the network. It is also shown in [7] that caches at the transmitters can improve load balancing and increase the opportunities for interference alignment. More recently, the authors in [8] evaluated the performance of cellular networks with edge caching via a hypergraph coloring problem. Furthermore, in [9], the authors studied the problem of maximizing the delivery rate of a fog radio access network for arbitrary prefetching strategies.

In this paper, we consider a general network setting with caches at both transmitters and receivers, and demonstrate how one can utilize caches at both transmitters and receivers to manage the interference and enhance the system performance in the physical layer. In particular, we consider a library of N files and a wireless network with K_T transmitters and K_R receivers, in which each transmitter and each receiver is equipped with a cache memory of a certain size. In particular, each transmitter and each receiver can cache up to M_T and M_R files, respectively. The system operates in two phases. The first phase is called the prefetching phase, where each cache is populated up to its limited size from the content of the library. This phase is followed by a delivery phase, where each user reveals its request for a file in the library. The transmitters then need to deliver the requested files to the receivers. Note that in the prefetching phase, the system is still unaware of the files that the receivers will request in the delivery phase. The goal is to design the cache contents in the prefetching phase and communication scheme

A short version of this paper will be presented at the IEEE International Symposium on Information Theory (ISIT), 2016. This work is the outcome of a collaboration that started while N. Naderializadeh was a research intern at Bell Labs. This work is in part supported by NSF grants CAREER 1408639, NETS-1419632, EARS-1411244, and ONR award N000141612189.

in the delivery phase to achieve the maximum throughput for arbitrary set of requested files. Due to their practical appeal, in this work we focus on one-shot linear delivery strategies. Interestingly, many of the previous works on caching have relied on one-shot schemes for content delivery (see, e.g. [4, 8]).

Our main result in this paper is the characterization of the one-shot linear sum degrees-of-freedom (sum-DoF) of the network, i.e., number of the receivers that can be served interference-free simultaneously, within a factor of 2 for all system parameters. In fact, we show that the one-shot linear sum-DoF of $\min \left\{ \frac{K_T M_T + K_R M_R}{N}, K_R \right\}$ is achievable, and this is within a factor of 2 of the optimum. This result shows that the one-shot linear sum-DoF of the network grows linearly with the aggregate cache size in the network (i.e., the cumulative memory available at all nodes). It also implies that caches at the transmitters' side are equally valuable as the caches on the receivers' side in the one-shot linear sum-DoF of the network. Our result, therefore, establishes a fundamental limit on the performance of one-shot delivery schemes for cache-aided interference management.

To achieve the aforementioned sum-DoF, we propose a particular pattern in cache placement so that each piece of each file in the library is available in the caches of $\frac{K_T M_T}{N}$ transmitters and $\frac{K_R M_R}{N}$ receivers. Once caching is done this way, we can show that for delivering any set of requested contents to the receivers, $\min \left\{ \frac{K_T M_T + K_R M_R}{N}, K_R \right\}$ of the receivers can be served at each time, interference-free. This gain is achieved by simultaneously exploiting the opportunity of collaborative interference cancellation (i.e. zero-forcing) at the transmitters' side and opportunity of eliminating known interference contributions at the receivers' side. The first opportunity is created by caching the pieces of each file at several transmitters. The second opportunity is available since pieces of a file requested by one user has been cached at some other receivers, and thus do not impose interference at those receivers effectively. Our proposed cache placement pattern maximizes the overall gain achieved by these opportunities for any arbitrary set of receiver requests and this gain can be achieved even with a simple one-shot linear delivery scheme.

Moreover, we demonstrate that our achievable sum-DoF is within a factor of 2 of the optimal sum-DoF for one-shot linear schemes. To prove the outer bound, we take a four-step approach in order to lower bound the number of communication blocks needed to deliver any set of requested files to the receivers. First, we show that the network can be converted to a virtual MISO interference channel in each block of communication. Using this conversion, we next write an integer optimization problem for the minimum number of communication blocks needed to deliver a fixed set of requests for a given caching realization. We then show how we can focus on average demands instead of the worst-case demands to derive an outer optimization problem on the number of communication blocks optimized over the caching realizations. Finally, we present a lower bound on the value of the aforementioned optimization problem, which leads to the desired upper bound on the one-shot linear sum-DoF of the network. This result illustrates that in this setting, caches at transmitters' side are equally valuable as caches at receivers' side. It also shows that caching offers a throughput gain that scales linearly with the size of the network.

The rest of the paper is organized as follows. We present the problem formulation in Section II. We state the main result in Section III. We prove the achievability of our main result in Section IV and the converse in Section V. Finally, we conclude the paper in Section VI.

II. PROBLEM FORMULATION

In this section, we first provide a high-level description of the problem setting and the main parameters in the system model, and then we present a detailed description of the problem formulation.

A. Problem Overview

Consider a wireless network, as illustrated in Figure 1, with K_T transmitters and K_R receivers, and also a library of N files, each of which contains F packets, where each packet is a vector of B bits. Each node in the network is equipped with a local cache memory of a certain size that can be used to cache contents arbitrarily from the library before the receivers reveal their requests and communication begins. In particular, each transmitter and each receiver is equipped with a cache of size $M_T F$ and $M_R F$ packets, respectively.

We assume that the system operates in two phases, namely the *prefetching phase* and the *delivery phase*. In the prefetching phase, each node can cache contents arbitrarily from the library subject to its cache size constraint. In particular, each transmitter selects up to $M_T F$ packets out of the entire library to store in its cache, and each receiver selects up to $M_R F$ packets out of the entire library to store in its cache. In the delivery phase, each receiver requests an arbitrary file from the library. Since each receiver may have cached parts of its desired file in the prefetching phase, the transmitters need to deliver the rest of the requested packets to the receivers over the wireless channel.

We assume that at each time, the transmitters employ a one-shot linear scheme, where a subset of requested packets are selected to be delivered interference-free to a corresponding subset of receivers. Each transmitter transmits a linear combination of the subset of the selected packets which it has cached in the prefetching phase. The interference is cancelled with the aid of cached contents as follows. Since each requested packet may be cached at multiple transmitters, the transmitters can collaborate in order to zero-force the outgoing interference of that packet at some of the unintended receivers. Moreover, the receivers can also use their cached packets as side information to eliminate the remaining incoming interference from to undesired packets.

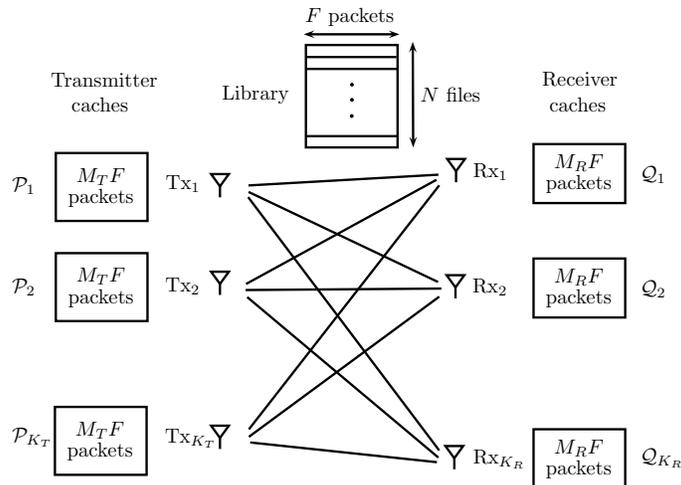


Fig. 1. Wireless network with K_T transmitters and K_R receivers, where each transmitter and each receiver caches up to $M_T F$ packets and $M_R F$ packets, respectively, from a library of N files, each composed of F packets.

Our objective is to design a cache placement scheme and a delivery scheme which maximize the number of packets that can be delivered at each time interference-free.

In this setting, we define the one-shot linear sum-degrees of freedom as the ratio of the number of delivered packets over the number of blocks needed for communicating those packets for any set of receiver demands. Finally, we define the one-shot linear sum-DoF of the network, denoted by $\text{DoF}_{\text{L,sum}}^*(N, M_T, M_R)$, as the maximum achievable one-shot linear sum-DoF over all caching realizations.

B. Detailed Problem Description

We consider a discrete-time additive white Gaussian noise channel, as illustrated in Figure 1, with K_T transmitters denoted by $\{\text{Tx}_i\}_{i=1}^{K_T}$ and K_R receivers denoted by $\{\text{Rx}_i\}_{i=1}^{K_R}$. The communication at time t over this channel is modeled by

$$Y_j(t) = \sum_{i=1}^{K_T} h_{ji} X_i(t) + Z_j(t), \quad (1)$$

where $X_i(t) \in \mathbb{C}$ denotes the signal transmitted by $\text{Tx}_i, i \in [K_T] \triangleq \{1, \dots, K_T\}$ and $Y_j(t)$ denotes the receive signal by $\text{Rx}_j, j \in [K_R]$. Moreover, $h_{ji} \in \mathbb{C}$ denotes the channel gain from Tx_i to Rx_j , assumed to stay fixed over the course of communication, and $Z_j(t)$ denotes the additive white Gaussian noise at Rx_j at time slot t , distributed as $\mathcal{CN}(0, 1)$. The transmit signal at $\text{Tx}_i, i \in [K_T] \triangleq \{1, \dots, K_T\}$ is subject to the power constraint $\mathbb{E}[|X_i(t)|^2] \leq P$.

We assume that each receiver will request an arbitrary file out of a library of N files $\{W_n\}_{n=1}^N$, which should be delivered by the transmitters. Each file W_n in the library contains F packets $\{\mathbf{w}_{n,f}\}_{f=1}^F$, where each packet is a vector of B bits; i.e., $\mathbf{w}_{n,f} \in \mathbb{F}_2^B$. Furthermore, we assume that each node in the network is equipped with a cache memory of a certain size that can be used to cache arbitrary contents from the library before the receivers reveal their requests and communication begins. In particular, each transmitter and each receiver is equipped with a cache of size $M_T F$ and $M_R F$ packets, respectively.

We assume that the network operates in two phases, namely the prefetching phase and the delivery phase, which are described in more detail as follows.

Prefetching Phase: In this phase, each node can store an arbitrary subset of the packets from the files in the library up to its cache size. In particular, each transmitter Tx_i chooses a subset \mathcal{P}_i of the NF packets in the library, where $|\mathcal{P}_i| \leq M_T F$, to store in its cache. Likewise, each receiver Rx_i stores a subset \mathcal{Q}_i of the packets in the library, where $|\mathcal{Q}_i| \leq M_R F$. Caching is done at the level of whole packets and we do not allow breaking the packets into smaller subpackets. Also, this phase takes place unaware of the receivers' future requests.

Delivery Phase: In this phase, each receiver $\text{Rx}_j, j \in [K_R]$, reveals its request for an arbitrary file W_{d_j} from the library for some $d_j \in [N]$. We let $\mathbf{d} = [d_1 \dots d_{K_R}]^T$ denote the demand vector. Depending on the demand vector \mathbf{d} and the cache contents, each receiver has already cached some packets of its desired file and there is no need to deliver them. The transmitters will be responsible for delivering the rest of the requested packets to the receivers. In order to make sure that any piece of content in the library is stored at the cache of at least one transmitter in the network, we assume that the transmitter cache size satisfies $K_T M_T \geq N$.

Each transmitter first employs a random Gaussian coding scheme $\psi: \mathbb{F}_2^B \rightarrow \mathbb{C}^{\tilde{B}}$ of rate $\log P + o(\log P)$ to encode each of its cached packets into a *coded packet* composed of \tilde{B} complex symbols, so that each coded packet carries one degree-of-freedom

(DoF). We denote the coded version of each packet $\mathbf{w}_{n,f}$ in the library by $\tilde{\mathbf{w}}_{n,f} \triangleq \psi(\mathbf{w}_{n,f})$. Afterwards, the communication takes place over H blocks, each of length \tilde{B} time slots. In each block $m \in [H]$, the goal is to deliver a subset of the requested packets, denoted by \mathcal{D}_m , to a subset of receivers, denoted by \mathcal{R}_m , such that each packet in \mathcal{D}_m is intended to exactly one of the receivers in \mathcal{R}_m . In addition, the set of transmitted packets in all blocks and the cache contents of the receivers should satisfy

$$\{\mathbf{w}_{d_j,f}\}_{f=1}^F \subset \left(\bigcup_{m=1}^H \mathcal{D}_m \right) \cup \mathcal{Q}_j, \quad \forall j \in [K_R], \quad (2)$$

which implies that for any receiver $\text{Rx}_j, j \in [K_R]$, each of its requested packets should be either transmitted in one of the blocks or already stored in its own cache.

In each block $m \in [H]$, we assume a one-shot linear scheme where each transmitter transmits an arbitrary linear combination of a subset of the coded packets in \mathcal{D}_m that it has cached. Particularly, $\text{Tx}_i, i \in [K_T]$ transmits $\mathbf{x}_i[m] \in \mathbb{C}^{\tilde{B}}$, where

$$\mathbf{x}_i[m] = \sum_{\substack{(n,f): \\ \mathbf{w}_{n,f} \in \mathcal{P}_i \cap \mathcal{D}_m}} v_{i,n,f}[m] \tilde{\mathbf{w}}_{n,f}, \quad (3)$$

and $v_{i,n,f}[m]$'s denote the complex beamforming coefficients that Tx_i uses to linearly combine its coded packets in block m .

On the receivers' side, the received signal of each receiver $\text{Rx}_j \in \mathcal{R}_m$ in block m , denoted by $\mathbf{y}_j[m] \in \mathbb{C}^{\tilde{B}}$, can be written as

$$\mathbf{y}_j[m] = \sum_{i=1}^{K_T} h_{ji} \mathbf{x}_i[m] + \mathbf{z}_j[m], \quad (4)$$

where $\mathbf{z}_j[m] \in \mathbb{C}^{\tilde{B}}$ denotes the noise vector at Rx_j in block m . Then, receiver Rx_j will use the contents of its cache to cancel (subtract out) the interference of some of undesired packets in \mathcal{D}_m , if they exist in its cache. In particular, each receiver $\text{Rx}_j \in \mathcal{R}_m$, forms a linear combination $\mathcal{L}_{j,m}$, as

$$\mathcal{L}_{j,m}(\mathbf{y}_j[m], \tilde{\mathcal{Q}}_j) \quad (5)$$

to recover $\tilde{\mathbf{w}}_{d_j,f} \in \mathcal{D}_m$, where $\tilde{\mathcal{Q}}_j$ denotes the set of coded packets cached at receiver Rx_j .

The communication in block $m \in H$ to transmit the packets in \mathcal{D}_m is successful, if there exist linear combinations (3) at the transmitters' side and (5) at receivers' side, such that for all $\text{Rx}_j \in \mathcal{R}_m$,

$$\mathcal{L}_{j,m}(\mathbf{y}_j[m], \tilde{\mathcal{Q}}_j) = \tilde{\mathbf{w}}_{d_j,f} + \mathbf{z}_j[m]. \quad (6)$$

The channel created in (6) is a point-to-point channel, whose capacity is $\log P + o(\log P)$. Hence, since each coded packet $\tilde{\mathbf{w}}_{d_j,f}$ is coded with rate $\log P + o(\log P)$, it can be decoded with vanishing error probability as B increases. We assume that the communication continues for H blocks until all the desired packets are successfully delivered to all receivers.

Since each packet carries one degree-of-freedom, the one-shot linear sum-degrees-of-freedom (sum-DoF) of $|\mathcal{D}_m|$ is achievable in each block $m \in [H]$. This implies that throughout the H blocks of communication, the one-shot linear sum-DoF of $\frac{|\bigcup_{m=1}^H \mathcal{D}_m|}{H}$ is achievable. Therefore, for a given caching realization, we define the one-shot linear sum-DoF to be maximum achievable one-shot linear sum-DoF for the worst case demands; i.e.,

$$\text{DoF}_{\text{L,sum}}^{\left(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}\right)} = \inf_{\mathbf{d}} \sup_{H, \{\mathcal{D}_m\}_{m=1}^H} \frac{\left| \bigcup_{m=1}^H \mathcal{D}_m \right|}{H}. \quad (7)$$

This leads us to the definition of the one-shot linear sum-DoF of the network as follows.

Definition 1. For a network with a library N files, each containing F packets, and cache size of M_T and M_R files at each transmitter and receiver, respectively, we define the one-shot linear sum-DoF of the network as the maximum achievable one-shot linear sum-DoF over all caching realizations; i.e.,

$$\text{DoF}_{\text{L,sum}}^*(N, M_T, M_R) = \sup_{\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}} \text{DoF}_{\text{L,sum}}^{\left(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}\right)} \quad (8)$$

$$\text{s.t. } |\mathcal{P}_i| \leq M_T F, \quad \forall i \in [K_T] \quad (9)$$

$$|\mathcal{Q}_i| \leq M_R F, \quad \forall i \in [K_R], \quad (10)$$

where $\text{DoF}_{\text{L,sum}}^{\left(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}\right)}$ is defined in (7).

III. MAIN RESULT AND ITS IMPLICATIONS

In this section, we present our main result on the one-shot linear sum-DoF of the network and its implications.

Theorem 1. *For a network with a library of N files, each containing F packets, and cache size of M_T and M_R files at each transmitter and each receiver, respectively, the one-shot linear sum-DoF of the network, as defined in Definition 1, satisfies*

$$\min \left\{ \frac{K_T M_T + K_R M_R}{N}, K_R \right\} \leq \text{DoF}_{\text{L,sum}}^*(N, M_T, M_R) \leq \min \left\{ 2 \frac{K_T M_T + K_R M_R}{N}, K_R \right\}, \quad (11)$$

for sufficiently large F .

In the following, we highlight the implications of Theorem 1 and its connections to some prior works:

- 1) (*Within a factor of 2 characterization*) The upper bound in (11) is within a factor of 2 of the lower bound in (11). Therefore, Theorem 1 characterizes the one-shot linear sum-DoF of a cache-aided wireless network to within a factor of 2, for all system parameters.
- 2) (*Aggregate cache size matters*) The one-shot linear sum-DoF characterized in Theorem 1 is proportional to the *aggregate* cache size that is available throughout the network, even-though these caches are *isolated*.
- 3) (*Equal contribution of transmitter and receiver caches*) Perhaps interestingly, the caches at both sides of the network, i.e., the transmitters' side and the receivers' side, are equally valuable in the achievable one-shot linear sum-DoF of the network. Note that in practice, size of each transmitter's cache, M_T , could be large. However, the number of transmitters (e.g., base stations) K_T is often small. On the other hand, size of the cache M_R at the receivers (e.g., cellphones) is small, whereas the number of receivers K_R is large. Therefore $K_T M_T$ could be comparable with $K_R M_R$. Our result in Theorem 1 shows that neither caches at the transmitters nor caches at the receivers should be ignored.
- 4) (*Linear scaling of DoF with network size*) Letting $K_T = K_R = K$, we observe that the one-shot linear sum-DoF scales *linearly* with the number of users in a *fully-connected* interference channel. Note that without caches, the one-shot linear sum-DoF of a fully-connected interference channel is bounded by 2, as shown in [14]. Hence, caching enables linear growth of the DoF without the need for more complex physical layer schemes.
- 5) (*Role of transmitter and receiver caches*) As we will show in Section IV, in (11), $\frac{K_T M_T}{N}$ represents the contribution of collaborative zero-forcing at the transmitters' side, and $\frac{K_R M_R}{N}$ represents the gain of canceling the known interference at the receivers' side.
- 6) (*Connection to single-server coded caching* [2]) A special case of our network model is the case with a single transmitter, which was previously considered in [2]. In this case, it can be shown that a sum-DoF of $\min \left\{ 1 + \frac{K_R M_R}{N}, K_R \right\}$ is achievable, which is equivalent to the global caching gain introduced in [2], indicating the number of receivers in the network that can be served simultaneously, interference-free. Hence, our result subsumes the result of [2] by generalizing it for the case of multiple transmitters.
- 7) (*Connection to multi-server coded caching* [3]) Another special case of our network model is the case where each transmitter has space to cache the entire library; i.e., $M_T = N$. This case was previously considered in [3] and it can be verified that in this case, a sum-DoF of $\min \left\{ K_T + \frac{K_R M_R}{N}, K_R \right\}$ is achievable. Hence, our result can also be viewed as a generalization of the result in [3] where the cache size of each transmitter may be arbitrarily smaller than the entire library size.

Remark 1. In practice, the files in the library have nonuniform demands and some of them are more popular than the rest. In this case, our algorithm can be used to cache and deliver the N most popular files. If a user requests one of the remaining less popular files, it can be directly served by a central base station. The parameter N can be tuned, based on the popularity pattern of the contents, in order to attain the best average performance.

Example 1. As an illustrative example, consider a cellular network with 5 base stations as transmitters, each with a 10 TB memory and 100 cellphones as receivers, each with a 32 GB memory. Moreover, consider a library of the 1000 most popular movie titles on Netflix, each with size of 5 GB. Then, Theorem 1 implies that at each time, around 11 cellphones can be served simultaneously interference-free, no matter what their demands are, in contrast to the naive time-sharing scheme, where at each time only 1 cellphone can be served. \square

The rest of the paper is devoted to the proof of Theorem 1. In particular, we illustrate the achievable scheme in Section IV and we present the converse argument in Section V.

IV. ACHIEVABLE SCHEME

In this section, we prove the achievability of Theorem 1 by presenting an achievable scheme which utilizes the caches at the transmitters and receivers efficiently to exploit the zero-forcing and interference cancellation opportunities at the transmitters' and receivers' sides, respectively. In particular, we introduce a prefetching strategy which maximizes the gains attained by the aforementioned opportunities in the delivery phase, no matter what the receiver demands are.

We first explain our achievable scheme through a simple, illustrative example and then proceed to mention our general achievable scheme.

A. Description of the Achievable Scheme via an Example

Consider a system with $K_T = 3$ transmitters and $K_R = 3$ receivers, where each transmitter has space to cache $M_T = 2$ files and each receiver has space to cache $M_R = 1$ file. The library has $N = 3$ files $W_1 = A$, $W_2 = B$, and $W_3 = C$, each consisting of F packets.

In the following, we will describe the prefetching and delivery phases in detail.

Prefetching Phase: In this phase, each file $W_n, n \in [3]$ in the library is broken into $\binom{3}{2} \binom{3}{1} = 9$ disjoint subfiles $W_{n,\mathcal{T},\mathcal{R}}$ for any $\mathcal{T} \subseteq [K_T] = [3]$ and $\mathcal{R} \subseteq [K_R] = [3]$ such that $|\mathcal{T}| = 2$ and $|\mathcal{R}| = 1$, where each subfile consists of $F/9$ packets. Each subfile $W_{n,\mathcal{T},\mathcal{R}}$ is then stored at the caches of the two transmitters in \mathcal{T} and the single receiver in \mathcal{R} . For example, file A is broken into 9 subfiles as follows:

$$A_{12,1}, A_{12,2}, A_{12,3}, A_{13,1}, A_{13,2}, A_{13,3}, A_{23,1}, A_{23,2}, A_{23,3},$$

where $A_{12,1}$ is stored at transmitters Tx_1 and Tx_2 as well as receiver Rx_1 , $A_{12,2}$ is stored at transmitters Tx_1 and Tx_2 as well as receiver Rx_2 , etc. We do the same partitioning for files B and C , as well.

It is easy to verify that each transmitter caches 6 subfiles of each file, hence the total size of its cached content is $3 \cdot (6 \cdot F/9) = 2F$ packets which satisfies its memory constraint. Also, each receiver caches 3 subfiles of each file and its total cached content has size $3 \cdot (3 \cdot F/9) = F$ packets, hence satisfying its memory constraint. Note that in this phase, we are unaware of receivers' future requests.

Delivery Phase: In this phase, each receiver reveals its request for a file in the library. Without loss of generality, assume that receivers Rx_1, Rx_2 and Rx_3 request files $W_{d_1} = A$, $W_{d_2} = B$ and $W_{d_3} = C$, respectively. Note that each receiver has already stored 3 subfiles of its desired file in its own cache, and therefore the transmitters need to deliver the 6 remaining subfiles of each requested file. In particular, the following 18 subfiles need to be delivered by the transmitters to the requesting receivers:

$$\begin{aligned} &A_{12,2}, A_{12,3}, A_{13,2}, A_{13,3}, A_{23,2}, A_{23,3} \text{ to receiver } \text{Rx}_1, \\ &B_{23,3}, B_{13,1}, B_{12,3}, B_{23,1}, B_{13,3}, B_{12,1} \text{ to receiver } \text{Rx}_2, \\ &C_{13,1}, C_{23,2}, C_{23,1}, C_{12,2}, C_{12,1}, C_{13,2} \text{ to receiver } \text{Rx}_3. \end{aligned} \tag{12}$$

We now show that we can break the 18 subfiles in (12) into 6 sets, each containing 3 subfiles, such that the subfiles in each set can be delivered simultaneously to the receivers, interference-free. Such a partitioning is illustrated through the 6 steps in Figure 2, where each step takes $\frac{F}{9}$ blocks. In each step, 3 subfiles are delivered to all the receivers simultaneously, while all the inter-user interference can be eliminated. For example, in the first step, as in Figure 2-(a), subfiles $A_{12,2}, B_{23,3}, C_{13,1}$ are respectively delivered to receivers Rx_1, Rx_2 , and Rx_3 at the same time. In Figure 3, we show in detail how the interference is cancelled in this step. The transmit signals of transmitters Tx_1, Tx_2 and Tx_3 can be respectively written as

$$\begin{aligned} X_1 &= -h_{32}\tilde{A}_{12,2} + h_{23}\tilde{C}_{13,1}, \\ X_2 &= h_{31}\tilde{A}_{12,2} - h_{13}\tilde{B}_{23,3}, \\ X_3 &= -h_{21}\tilde{C}_{13,1} + h_{12}\tilde{B}_{23,3}, \end{aligned}$$

where for any subfile $W_{n,\mathcal{T},\mathcal{R}}, \tilde{W}_{n,\mathcal{T},\mathcal{R}}$ denotes its coded version. For simplicity, in this example, we ignore the power constraint at the transmitters. On the other hand, the received signals by receivers Rx_1, Rx_2 and Rx_3 can be respectively written as

$$\begin{aligned} Y_1 &= (h_{12}h_{31} - h_{11}h_{32})\tilde{A}_{12,2} + (h_{11}h_{23} - h_{13}h_{21})\tilde{C}_{13,1} + Z_1, \\ Y_2 &= (h_{23}h_{12} - h_{22}h_{13})\tilde{B}_{23,3} + (h_{22}h_{31} - h_{21}h_{32})\tilde{A}_{12,2} + Z_2, \\ Y_3 &= (h_{31}h_{23} - h_{33}h_{21})\tilde{C}_{13,1} + (h_{33}h_{12} - h_{32}h_{13})\tilde{B}_{23,3} + Z_3. \end{aligned}$$

Now, note that receivers Rx_1, Rx_2 and Rx_3 can cancel the interference due to $C_{13,1}, A_{12,2}$, and $B_{23,3}$, respectively, since they already have each respective subfile in their own cache. Therefore, all the interference in the network can be effectively eliminated and the receivers will be able to decode their desired subfiles. Likewise, one can verify that all the receivers can receive their desired subfiles interference-free in all the 6 steps of communication depicted in Figure 2.

Consequently, the 18 subfiles in (12), each of which consists of $F/9$ packets, are delivered to the receivers in 6 steps, each consisting of $F/9$ blocks. Note that our particular file splitting pattern in the prefetching phase and the particular scheduling pattern in the delivery phase allows us to maximally exploit the two gains of zero-forcing the outgoing interference on the transmitters' side and canceling the known interference on the receivers' side, no matter what the receiver demands are in the delivery phase. Therefore, the sum-DoF of $\frac{18 \cdot F/9}{6 \cdot F/9} = 3 = \min \left\{ \frac{K_T M_T + K_R M_R}{N}, K_R \right\}$ is achievable in this network.

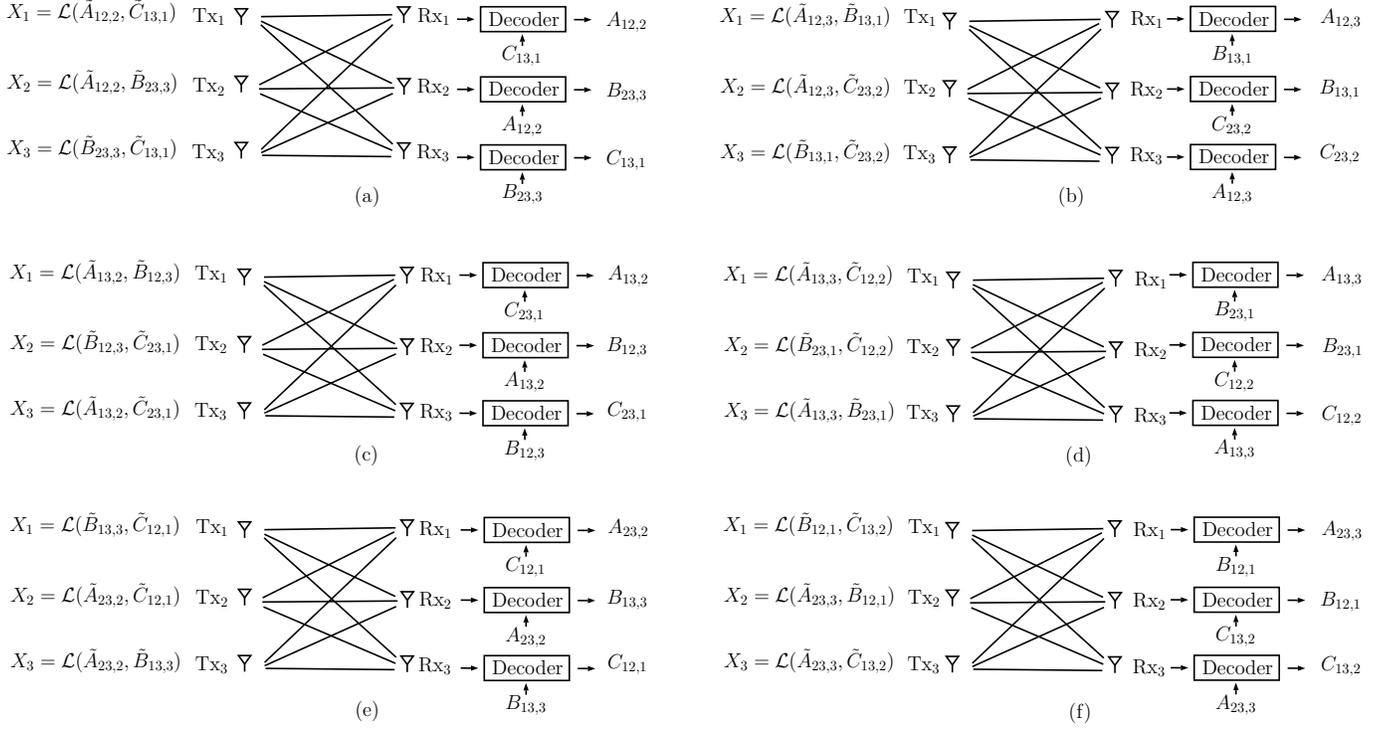


Fig. 2. Delivery phase for the example in Section IV-A for respective requests of files A , B and C by receivers Rx_1 , Rx_2 and Rx_3 , where $\mathcal{L}(\alpha, \beta)$ denotes some linear combination of α and β . In every step, each pair of transmitters collaborate to zero-force the interference due to a specific subfile at a certain undesired receiver. Moreover, each receiver also uses its cache contents to cancel the interference due to the other interfering packet. Therefore, the communication is interference-free in all 6 steps.

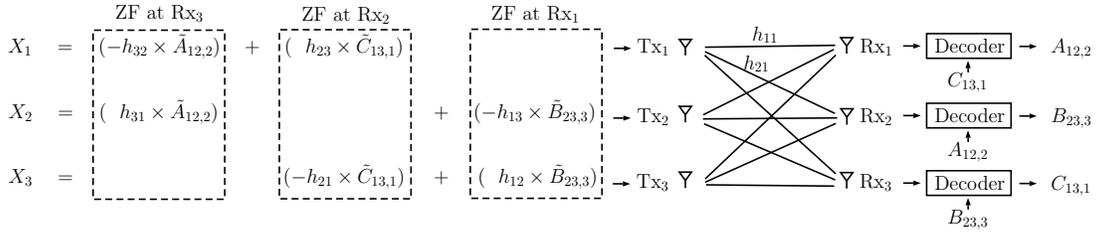


Fig. 3. More detailed description of the linear encoding and decoding schemes used in the delivery phase step in Figure 2-(a). In this step, Tx₁ and Tx₂ zero-force $A_{12,2}$ at Rx₃, Tx₁ and Tx₃ zero-force $C_{13,1}$ at Rx₂, and Tx₂ and Tx₃ zero-force $B_{23,3}$ at Rx₁. Moreover, Rx₁, Rx₂ and Rx₃ can cancel the interference due to $C_{13,1}$, $A_{12,2}$, and $B_{23,3}$, respectively, since they already have each respective subfile in their own cache.

B. Description of the General Achievable Scheme

Our general achievable scheme is given in Algorithm 1. In this algorithm, we use the notation

$$t_T \triangleq \frac{K_T M_T}{N}, \quad t_R \triangleq \frac{K_R M_R}{N}, \quad (13)$$

and for now, we assume that t_T and t_R are integers. Recall that in the example in Section IV-A, $t_T = 2$ and $t_R = 1$.

In the following, we will describe the prefetching and delivery phases in more detail.

1) *Prefetching Phase*: For any file W_n in the library, $n \in [N]$, we partition it into $\binom{K_T}{t_T} \binom{K_R}{t_R}$ disjoint subfiles of equal sizes¹, denoted by

$$W_n = \left\{ W_{n, \mathcal{T}, \mathcal{R}} \right\}_{\substack{\mathcal{T} \subseteq [K_T]: |\mathcal{T}|=t_T \\ \mathcal{R} \subseteq [K_R]: |\mathcal{R}|=t_R}} \quad (14)$$

Based on the above partitioning, in the prefetching phase, each transmitter Tx_{*i*} stores a subset \mathcal{P}_i of the packets in the

¹Due to the assumption that F is sufficiently large, we can assume that it is an integer multiple of $\binom{K_T}{t_T} \binom{K_R}{t_R}$.

Algorithm 1 Achievable scheme for Theorem 1

Prefetching Phase:

- 1: **for** $n = 1, \dots, N$
- 2: Partition W_n into $\binom{K_T}{t_T} \binom{K_R}{t_R}$ disjoint subfiles $\{W_{n,\mathcal{T},\mathcal{R}}\}_{\mathcal{T} \subseteq [K_T], |\mathcal{T}|=t_T, \mathcal{R} \subseteq [K_R], |\mathcal{R}|=t_R}$ of equal sizes.
- 3: **end**
- 4: **for** $i = 1, \dots, K_T$
- 5: Tx_{*i*} caches all $W_{n,\mathcal{T},\mathcal{R}}$ for which $i \in \mathcal{T}$.
- 6: **end**
- 7: **for** $j = 1, \dots, K_R$
- 8: Rx_{*j*} caches all $W_{n,\mathcal{T},\mathcal{R}}$ for which $j \in \mathcal{R}$.
- 9: **end**

Delivery Phase:

- 10: **for** $j \in [K_R]$
 - 11: **for** $\mathcal{T} \subseteq [K_T]$ s.t. $|\mathcal{T}| = t_T$
 - 12: **for** $\mathcal{R} \subseteq [K_R] \setminus \{j\}$ s.t. $|\mathcal{R}| = t_R$
 - 13: partition $W_{d_j,\mathcal{T},\mathcal{R}}$ to $\frac{t_R! [K_R - (t_R + 1)]!}{[K_R - (t_R + t_T)]!}$ disjoint subfiles $\left\{ W_{d_j,\mathcal{T},\pi,\pi'} \right\}_{\substack{\pi \in \Pi_{\mathcal{R}} \\ \pi' \in \Pi_{[K_R] \setminus (\mathcal{R} \cup \{j\}), t_T - 1}}}$ of equal sizes.
 - 14: **end**
 - 15: **end**
 - 16: **end**
 - 17: **for** $\mathcal{T} \subseteq [K_T]$ s.t. $|\mathcal{T}| = t_T$
 - 18: **for** $\mathcal{R} \subseteq [K_R]$ s.t. $|\mathcal{R}| = t_T + t_R$
 - 19: **for** $\pi \in \Pi_{\mathcal{R}}^{\text{circ}}$
 - 20: Each transmitter Tx_{*i*} transmits a linear combination of the coded subfiles as in $X_i = \mathcal{L}_{i,\mathcal{T},\pi} \left(\left\{ \tilde{W}_{d_{\pi(l)},\mathcal{T} \oplus_{K_T} (l-1), \pi[l+1:l+t_R], \pi[l+t_R+1:l+t_R+t_T-1]} : l \in [t_T + t_R], i \in \mathcal{T} \oplus_{K_T} (l-1) \right\} \right)$ using the linear combinations shown in Lemma 2 such that the subfiles $\left\{ W_{d_{\pi(l)},\mathcal{T} \oplus_{K_T} (l-1), \pi[l+1:l+t_R], \pi[l+t_R+1:l+t_R+t_T-1]} : l \in [t_T + t_R] \right\}$ are simultaneously delivered to the receivers in \mathcal{R} interference-free.
 - 21: **end**
 - 22: **end**
 - 23: **end**
-

library as described below.

$$\mathcal{P}_i = \{W_{n,\mathcal{T},\mathcal{R}} : i \in \mathcal{T}\}. \quad (15)$$

Illustration. For instance, in the example network considered in Section IV-A, transmitter Tx₃ stores the following subset of packets in its cache.

$$\begin{aligned} \mathcal{P}_3 &= \{W_{1,13,1}, W_{1,13,2}, W_{1,13,3}, W_{1,23,1}, W_{1,23,2}, W_{1,23,3}, \\ &\quad W_{2,13,1}, W_{2,13,2}, W_{2,13,3}, W_{2,23,1}, W_{2,23,2}, W_{2,23,3}, \\ &\quad W_{3,13,1}, W_{3,13,2}, W_{3,13,3}, W_{3,23,1}, W_{3,23,2}, W_{3,23,3}\} \\ &= \{A_{13,1}, A_{13,2}, A_{13,3}, A_{23,1}, A_{23,2}, A_{23,3}, \\ &\quad B_{13,1}, B_{13,2}, B_{13,3}, B_{23,1}, B_{23,2}, B_{23,3}, \\ &\quad C_{13,1}, C_{13,2}, C_{13,3}, C_{23,1}, C_{23,2}, C_{23,3}\}. \end{aligned} \quad \square$$

Based on the above caching strategy, we can verify that the total number of packets cached by transmitter Tx_{*i*} equals

$$N \binom{K_T - 1}{t_T - 1} \binom{K_R}{t_R} \frac{F}{\binom{K_T}{t_T} \binom{K_R}{t_R}} = NF \frac{t_T}{K_T} = M_T F \text{ packets,}$$

hence satisfying its memory size constraint, where $\binom{K_T - 1}{t_T - 1}$ is the number of subsets $\mathcal{T} \subseteq [K_T]$ of size t_T which include the transmitter index i .

Likewise, in the prefetching phase, each receiver Rx_{*j*} stores a subset \mathcal{Q}_j of the packets in the library as described below.

$$\mathcal{Q}_j = \{W_{n,\mathcal{T},\mathcal{R}} : j \in \mathcal{R}\}. \quad (16)$$

Illustration. For instance, in the example network considered in Section IV-A, receiver Rx₂ stores the following subset of

packets in its cache.

$$\begin{aligned} \mathcal{Q}_2 &= \{W_{1,12,2}, W_{1,13,2}, W_{1,23,2}, W_{2,12,2}, W_{2,13,2}, W_{2,23,2}, W_{3,12,2}, W_{3,13,2}, W_{3,23,2}\} \\ &= \{A_{12,2}, A_{13,2}, A_{23,2}, B_{12,2}, B_{13,2}, B_{23,2}, C_{12,2}, C_{13,2}, C_{23,2}\}. \end{aligned} \quad \square$$

This suggests that the total number of packets cached by receiver Rx_j is equal to

$$N \binom{K_T}{t_T} \binom{K_R - 1}{t_R - 1} \frac{F}{\binom{K_T}{t_T} \binom{K_R}{t_R}} = NF \frac{t_R}{K_R} = M_R F \text{ packets,}$$

which also satisfies its memory size constraint.

2) *Delivery Phase:* In this section, we first describe the delivery phase for the case where $t_T + t_R \leq K_R$, so that the first term in the lower bound in (11) is dominant. We will later show how to deal with the case where $t_T + t_R > K_R$.

In the delivery phase, the receiver requests are revealed, and in particular, each receiver $Rx_j, j \in [K_R]$ requests a file W_{d_j} from the library and the transmitters need to deliver the subfiles in

$$\{W_{d_j, \mathcal{T}, \mathcal{R}} : j \notin \mathcal{R}\}$$

to receiver Rx_j ; i.e., the subfiles of file W_{d_j} which have not been already stored in the cache of receiver Rx_j .

In the following, our goal is to show that the set of packets which need to be delivered to the receivers can be partitioned into subsets of size $t_T + t_R$ such that the packets in each subset can be scheduled together. To this end, we need to further break each subfile to smaller subfiles. In particular, for any $j \in [K_R], \mathcal{T} \subseteq [K_T]$ s.t. $|\mathcal{T}| = t_T, \mathcal{R} \subseteq [K_R] \setminus \{j\}$ s.t. $|\mathcal{R}| = t_R$, we partition $W_{d_j, \mathcal{T}, \mathcal{R}}$ to $\frac{t_R! [K_R - (t_R + 1)]!}{[K_R - (t_R + t_T)]!}$ smaller disjoint subfiles of equal sizes denoted by

$$W_{d_j, \mathcal{T}, \mathcal{R}} = \left\{ W_{d_j, \mathcal{T}, \pi, \pi'} \right\}_{\substack{\pi \in \Pi_{\mathcal{R}} \\ \pi' \in \Pi_{[K_R] \setminus (\mathcal{R} \cup \{j\}), t_T - 1}}}, \quad (17)$$

where for a set \mathcal{S} , $\Pi_{\mathcal{S}}$ denotes the set of permutations of \mathcal{S} , and for any $t \in \{1, \dots, |\mathcal{S}|\}$, $\Pi_{\mathcal{S}, t}$ denotes the set of all permutations of all subsets of \mathcal{S} of size t ; i.e.,

$$\Pi_{\mathcal{S}, t} = \bigcup_{\mathcal{A} \subseteq \mathcal{S}, |\mathcal{A}|=t} \Pi_{\mathcal{A}}.$$

Remark 2. Note that in the example setting discussed in Section IV-A, $\frac{t_R! [K_R - (t_R + 1)]!}{[K_R - (t_R + t_T)]!} = 1$, which implies that further partitioning of the subfiles is not needed.

The advantage of further breakdown of the subfiles in (17) is that we can now partition the set of the subfiles which need to be delivered to the receivers into certain subsets of size $t_T + t_R$ such that each subfile $W_{d_j, \mathcal{T}, \pi, \pi'}$ intended for receiver Rx_j is zero-forced at the receivers with indices in π' . Moreover, since this subfile is also already cached at the receivers with indices in π , the communication will be interference-free for each set of the $t_T + t_R$ subfiles.

We show how to do such a partitioning in Lemma 1. In this lemma, we use the following notation: For a set \mathcal{R} , we let $\Pi_{\mathcal{R}}^{\text{circ}}$ denote the set of $(|\mathcal{R}| - 1)!$ circular permutations of \mathcal{R} .² Moreover, for a set \mathcal{S} , a permutation $\pi \in \Pi_{\mathcal{S}}$ and two integers i, j satisfying $j \geq i$, we define $\pi[i : j]$ as

$$\pi[i : j] = [\pi(i \oplus_{|\mathcal{S}|} 0) \ \pi(i \oplus_{|\mathcal{S}|} 1) \ \pi(i \oplus_{|\mathcal{S}|} 2) \ \dots \ \pi(i \oplus_{|\mathcal{S}|} (j - i))],$$

where for an integer m , $i \oplus_m j$ is defined as

$$i \oplus_m j = 1 + (i + j - 1 \pmod{m}). \quad (18)$$

Finally, for a set \mathcal{T} and an integer j , we let $\mathcal{T} \oplus_m j$ denote entry-wise addition of elements of \mathcal{T} with j modulo m , as defined in (18).

Lemma 1. *Given the prefetching phase in Section IV-B1, for any receivers' demand vector \mathbf{d} , the set of subfiles which need to be delivered to the receivers can be partitioned into disjoint subsets of size $t_T + t_R$ as*

$$\bigcup_{\substack{\mathcal{T} \subseteq [K_T]: |\mathcal{T}|=t_T \\ \mathcal{R} \subseteq [K_R]: |\mathcal{R}|=t_T+t_R \\ \pi \in \Pi_{\mathcal{R}}^{\text{circ}}}} \left\{ W_{d_{\pi(l)}, \mathcal{T} \oplus_{K_T} (l-1), \pi[l+1:l+t_R], \pi[l+t_R+1:l+t_R+t_T-1]} : l \in [t_T + t_R] \right\}. \quad (19)$$

Proof. See Appendix A. □

²A circular permutation of a set \mathcal{R} is a way of arranging the elements of \mathcal{R} around a fixed circle. The number of distinct circular permutations of a set \mathcal{R} is equal to $(|\mathcal{R}| - 1)!$. For example, if $\mathcal{R} = \{1, 2, 3\}$, then $\Pi_{\mathcal{R}}^{\text{circ}} = \{[1, 2, 3], [1, 3, 2]\}$.

Illustration. For the example network mentioned in Section IV-A, the set of 18 subfiles which need to be delivered to the receivers, as in (12), can be partitioned to the following 6 sets.

$$\begin{aligned} & \{A_{12,2}, B_{23,3}, C_{13,1}\} \cup \{A_{12,3}, B_{13,1}, C_{23,2}\} \cup \{A_{13,2}, B_{12,3}, C_{23,1}\} \\ & \cup \{A_{13,3}, B_{23,1}, C_{12,2}\} \cup \{A_{23,2}, B_{13,3}, C_{12,1}\} \cup \{A_{23,3}, B_{12,1}, C_{13,2}\}. \end{aligned} \quad (20)$$

Based on the partitioning of the small subfiles that need to be delivered to the receivers in Lemma 1, we will have $\binom{K_T}{t_T} \binom{K_R}{t_T+t_R} (t_T+t_R-1)!$ steps of communication, where at each step, specific sets \mathcal{T} and \mathcal{R} and a permutation π are fixed as in (19), and each transmitter Tx_i will transmit a linear combination of the coded subfiles for which $i \in \mathcal{T} \oplus_{K_T} (l-1)$; i.e.,

$$X_i = \mathcal{L}_{i,\mathcal{T},\pi} \left(\left\{ \tilde{W}_{d_{\pi(l), \mathcal{T} \oplus_{K_T} (l-1), \pi[l+1:l+t_R], \pi[l+t_R+1:l+t_R+t_T-1]}} : l \in [t_T+t_R], i \in \mathcal{T} \oplus_{K_T} (l-1) \right\} \right), \quad (21)$$

where for any subfile $W_{d_j, \mathcal{T}, \pi, \pi'}$, $\tilde{W}_{d_j, \mathcal{T}, \pi, \pi'}$ denotes the corresponding coded subfile containing PHY coded symbols, and $\mathcal{L}_{i,\mathcal{T},\pi}(\cdot)$ represents the linear combination that transmitter Tx_i chooses for sending the subfiles in (21).

We will next show that under such a delivery scheme, there always exists a choice of linear combinations at the transmitters so that at each step, the communication will be interference-free and all the t_T+t_R receivers in \mathcal{R} can decode their desired packets, as we also showed in the example setting in Section IV-A.

Lemma 2. *For any subset of t_T transmitters $\mathcal{T} \subseteq [K_T]$, any subset of t_T+t_R receivers $\mathcal{R} \subseteq [K_R]$, and any circular permutation $\pi \in \Pi_{\mathcal{R}}^{\text{circ}}$, there exists a choice of the linear combinations $\{\mathcal{L}_{i,\mathcal{T},\pi}(\cdot)\}_{i=1}^{K_T}$ in (21) such that the set of t_T+t_R subfiles in*

$$\left\{ W_{d_{\pi(l), \mathcal{T} \oplus_{K_T} (l-1), \pi[l+1:l+t_R], \pi[l+t_R+1:l+t_R+t_T-1]}} : l \in [t_T+t_R] \right\}, \quad (22)$$

can be delivered simultaneously and interference-free by the transmitters in $\bigcup_{l \in [t_T+t_R]} (\mathcal{T} \oplus_{K_T} (l-1))$ to the receivers in \mathcal{R} .

Proof. For ease of notation and without loss of generality, assume

$$\mathcal{T} = \{1, \dots, t_T\}, \mathcal{T} \oplus_{K_T} (l-1) = \{l, \dots, t_T+l\}, \mathcal{R} = \{1, \dots, t_T+t_R\}, \pi = [1, \dots, t_T+t_R].$$

First, we need to determine the subset of the subfiles which is available at each transmitter. It is easy to verify that

- If $i \in \{1, \dots, t_T-1\}$, then transmitter Tx_i has subfiles

$$\left\{ \tilde{W}_{d_{\pi(l), \mathcal{T} \oplus_{K_T} (l-1), \pi[l+1:l+t_R], \pi[l+t_R+1:l+t_R+t_T-1]}} : l \in \{1, \dots, i\} \right\}; \quad (23)$$

- If $i \in \{t_T, \dots, t_T+t_R\}$, then transmitter Tx_i has subfiles

$$\left\{ \tilde{W}_{d_{\pi(l), \mathcal{T} \oplus_{K_T} (l-1), \pi[l+1:l+t_R], \pi[l+t_R+1:l+t_R+t_T-1]}} : l \in \{i-t_T+1, \dots, i\} \right\}; \quad (24)$$

- and if $i \in \{t_T+t_R+1, \dots, 2t_T+t_R-1\}$, then transmitter Tx_i has subfiles

$$\left\{ \tilde{W}_{d_{\pi(l), \mathcal{T} \oplus_{K_T} (l-1), \pi[l+1:l+t_R], \pi[l+t_R+1:l+t_R+t_T-1]}} : l \in \{i-t_T+1, \dots, t_T+t_R\} \right\}. \quad (25)$$

Since each transmitter sends a linear combination of the subfiles that it has, the transmit signal of transmitter Tx_i can be written as

$$X_i = \begin{cases} \sum_{l=1}^i v_{i,l} \tilde{W}_{d_l, \{1, \dots, t_T+l\}, \{l+1, \dots, l+t_R\}, \{l+t_R+1, \dots, l+t_R+t_T-1\}}, & \text{if } i \in \{1, \dots, t_T-1\} \\ \sum_{l=i-t_T+1}^i v_{i,l} \tilde{W}_{d_l, \{1, \dots, t_T+l\}, \{l+1, \dots, l+t_R\}, \{l+t_R+1, \dots, l+t_R+t_T-1\}}, & \text{if } i \in \{t_T, \dots, t_T+t_R\} \\ \sum_{l=i-t_T+1}^{t_T+t_R} v_{i,l} \tilde{W}_{d_l, \{1, \dots, t_T+l\}, \{l+1, \dots, l+t_R\}, \{l+t_R+1, \dots, l+t_R+t_T-1\}}, & \text{if } i \in \{t_T+t_R+1, \dots, 2t_T+t_R-1\} \end{cases}. \quad (26)$$

This implies that the received signal at receiver Rx_j , $j \in \{1, \dots, t_T+t_R\}$ can be written as

$$\begin{aligned} Y_j &= \sum_{i=1}^{2t_T+t_R-1} h_{ji} X_i + Z_j \\ &= \sum_{i=j}^{t_T+j} h_{ji} v_{i,j} \tilde{W}_{d_j, \{j, \dots, t_T+j\}, \{j+1, \dots, j+t_R\}, \{j+t_R+1, \dots, j+t_R+t_T-1\}} \\ &\quad + \sum_{l=j+1}^{j+t_T-1} \sum_{i=l}^{t_T+l} h_{ji} v_{i,l} \tilde{W}_{d_l, \{1, \dots, t_T+l\}, \{l+1, \dots, l+t_R\}, \{l+t_R+1, \dots, l+t_R+t_T-1\}} \end{aligned} \quad (27)$$

$$+ \sum_{l=j-t_R}^{j-1} \sum_{i=l}^{t_T+l} h_{ji} v_{i,l} \tilde{W}_{d_i, \{l, \dots, t_T+l\}, \{l+1, \dots, l+t_R\}, \{l+t_R+1, \dots, l+t_R+t_T-1\}} + Z_j. \quad (28)$$

Now, note that in (28), the first term corresponds to the desired subfile of receiver Rx_j , while the second and third terms correspond to the undesired subfiles whose interference needs to be canceled at this receiver. However, note that the subfiles in the third term are already cached at receiver Rx_j and hence it is able to cancel their incoming interference. Hence, in order for all receivers $Rx_j, j \in \{1, \dots, t_T + t_R\}$ to receive their subfiles interference-free, there should exist a choice of linear combination coefficients $\{v_{i,l}\}$ such that

$$\sum_{i=j}^{t_T+j} h_{ji} v_{i,j} = 1, \forall j \in \{1, \dots, t_T + t_R\} \quad (29)$$

$$\sum_{i=l}^{t_T+l} h_{ji} v_{i,l} = 0, \forall j \in \{1, \dots, t_T + t_R\}, \forall l \in \{j+1, \dots, j+t_T-1\}. \quad (30)$$

Equations (29)-(30) introduce a system of $t_T(t_T + t_R)$ linear equations. On the other hand, the number of variables $\{v_{i,l}\}$ is also equal to $t_T(t_T + t_R)$. This indicates that there always exists a choice of linear combination coefficients $\{v_{i,l}\}$ such that (29)-(30) are satisfied. Finally, note that by scaling all the transmit signals by a large enough factor, the power constraint at all the transmitters can also be satisfied. Hence the proof is complete. \square

Remark 3. As mentioned in Section II, we assume that the channel gains remain constant over the course of communication. However, for the delivery scheme presented in the proof of Lemma 2, this assumption can be relaxed, since we only need the channel gains to remain unchanged for each block of communication and they can be allowed to vary among different blocks.

Remark 4. In the delivery scheme presented in the proof of Lemma 2, we only used zero-forcing at the transmitters in order to cancel their outgoing interference, which is DoF-optimal. However, in general one can use any scheme that exploits the collaboration among the transmitters in order to optimize the actual rates in the finite-SNR regime (such as the schemes suited for the MIMO broadcast channels [15]).

C. Analysis of the Sum-DoF of the Proposed Achievable Scheme

As a result of Lemmas 1 and 2, it is clear that for any set of receiver demands in the delivery phase, we can schedule all the requested subfiles in groups of size $t_T + t_R$. Now, if t_T and/or t_R are not integers, we can split the memories and the files proportionally so that for each new partition, the aforementioned scheme can be applied for updated t_T and t_R which are integers. Hence, combining the schemes over different partitions allows us to serve $t_T + t_R$ simultaneously, interference-free, for any values of t_T and t_R such that $t_T + t_R \leq K_R$.³

Finally, if $t_T + t_R > K_R$, then since we cannot serve more than K_R receivers, we can neglect some of the caches at either the transmitters' side or the receivers' side and use a fraction of the caches with new sizes $\frac{N}{K_T} \leq M'_T \leq M_T$ and $M'_R \leq M_R$ so that $\frac{K_T M'_T + K_R M'_R}{N} = K_R$. We can then use Algorithm 1 to serve all the K_R receivers simultaneously without interference.

As we showed in Section IV-B1, our prefetching phase respects the cache size constraint of all the transmitters and receivers. Moreover, given our prefetching phase, each receiver Rx_j caches $\binom{K_T}{t_T} \binom{K_R-1}{t_R-1} \frac{F}{\binom{K_T}{t_T} \binom{K_R}{t_R}} = \frac{M_R}{N} F$ packets of each file in the library. Hence, for each set of requested files by the receivers, a total of $K_R \left(1 - \frac{M_R}{N}\right) F$ packets need to be delivered by the transmitters to the receivers.

Therefore, based on the delivery phase mentioned in Section IV-B, the number of blocks required to deliver all the $K_R \left(1 - \frac{M_R}{N}\right) F$ packets to the receivers is equal to $\frac{K_R \left(1 - \frac{M_R}{N}\right) F}{\min\{t_T+t_R, K_R\}}$. This suggests that for any set of receiver demands, sum-DoF of

$$\frac{K_R \left(1 - \frac{M_R}{N}\right) F}{\min\{t_T+t_R, K_R\}} = \min\{t_T + t_R, K_R\} = \min \left\{ \frac{K_T M_T + K_R M_R}{N}, K_R \right\}$$

is achievable, hence completing the proof of achievability of Theorem 1.

V. CONVERSE

In this section, we prove the converse of Theorem 1. In particular, we show that the lower bound on the one-shot linear sum-DoF in (11) is within a factor of 2 of the optimal one-shot linear sum-DoF. In order to prove the converse, we take four steps as detailed in the following sections. First, we demonstrate how in each block of communication, the network can be converted into a virtual MISO interference channel. Second, we use this conversion to write an integer optimization problem

³In [2], this method is referred to as *memory-sharing*, which resembles time-sharing in network information theory.

for the minimum number of communication blocks needed to deliver a set of receiver demands for a given caching realization. Third, we show how we can focus on average demands instead of the worst-case demands to derive an outer optimization problem on the number of communication blocks optimized over the caching realizations. Finally, we present a lower bound on the value of the aforementioned outer optimization problem, which leads to the desired upper bound on the one-shot linear sum-DoF of the network.

A. Conversion to a Virtual MISO Interference Channel

Consider any caching realization $(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R})$ and any demand vector \mathbf{d} . As discussed in Section II, in each communication block a subset of requested packets are selected to be sent to a corresponding subset of distinct receivers. Now, we can state the following lemma, which bounds the number of packets that can be scheduled together in a single communication block using a one-shot linear scheme.

Lemma 3. *Consider a single communication block where a set $\{\mathbf{w}_{n_l, f_l}\}_{l=1}^L$ of L packets are scheduled to be transmitted together to L distinct receivers. In order for each receiver to successfully decode its desired packet, the number of these concurrently-scheduled packets should be bounded by*

$$L \leq \min_{l \in [L]} |\mathcal{T}_l| + |\mathcal{R}_l|, \quad (31)$$

where for any $l \in [L]$, \mathcal{T}_l and \mathcal{R}_l denote the set of transmitters and receivers which have cached the packet \mathbf{w}_{n_l, f_l} , respectively.

Proof. For ease of notation and without loss of generality, suppose that in the considered block, L packets $\{\mathbf{w}_{1,1}, \dots, \mathbf{w}_{L,1}\}$ are scheduled to be sent to L receivers $\{\mathbf{R}x_1, \dots, \mathbf{R}x_L\}$, respectively. Each transmitter $\mathbf{T}x_i, i \in [K_T]$ will transmit

$$\mathbf{x}_i = \sum_{l: i \in \mathcal{T}_l} v_{i,l,1} \tilde{\mathbf{w}}_{l,1}, \quad (32)$$

where we have dropped the dependency on the block index, since we are focusing on a single block. On the other hand, the received signal of receiver $\mathbf{R}x_j, j \in [L]$ can be written as

$$\mathbf{y}_j = \sum_{i=1}^{K_T} h_{ji} \mathbf{x}_i \quad (33)$$

$$= \sum_{i=1}^{K_T} h_{ji} \sum_{l: i \in \mathcal{T}_l} v_{i,l,1} \tilde{\mathbf{w}}_{l,1} \quad (34)$$

$$= \sum_{l=1}^L \sum_{i \in \mathcal{T}_l} h_{ji} v_{i,l,1} \tilde{\mathbf{w}}_{l,1}. \quad (35)$$

Therefore, (35) implies that we can effectively convert the network into a new MISO interference channel with L virtual transmitters $\{\widehat{\mathbf{T}}x_l\}_{l=1}^L$, where $\widehat{\mathbf{T}}x_l$ is equipped with $|\mathcal{T}_l|$ antennas, and L single-antenna receivers $\{\mathbf{R}x_j\}_{j=1}^L$, in which each virtual transmitter $\widehat{\mathbf{T}}x_l$ intends to send the coded packet $\tilde{\mathbf{w}}_{l,1}$ to receiver $\mathbf{R}x_l$. Each antenna in the new network corresponds to a transmitter in the original network. Hence, the channel vectors are correlated in the new network. In fact, as (35) suggests, all the antennas corresponding to the same transmitter in the original network have the same channel gain vectors to the receivers in the new network.

In the constructed MISO interference channel, we take a similar approach as in [14] in order to bound the one-shot linear sum-DoF of the network. Each virtual transmitter $\widehat{\mathbf{T}}x_l$ in the constructed MISO network will select a beamforming vector $\mathbf{v}_l \in \mathbb{C}^{|\mathcal{T}_l| \times 1}$ (which consists of the coefficients chosen by the original transmitters corresponding to its antennas) to transmit its desired symbol. Denoting the channel gain vector between transmitter $\widehat{\mathbf{T}}x_l$ and receiver $\mathbf{R}x_j$ as $\mathbf{h}_{jl} \in \mathbb{C}^{|\mathcal{T}_l| \times 1}$, the decodability conditions can be written as

$$\mathbf{h}_{jl}^T \mathbf{v}_l = 0, \quad \forall l \neq j \text{ s.t. } j \notin \mathcal{R}_l \quad (36)$$

$$\mathbf{h}_{jj}^T \mathbf{v}_j \neq 0, \quad \forall j \in [L]. \quad (37)$$

Now, each of the vectors $\mathbf{v}_l, l \in [L]$ can be written as

$$\mathbf{v}_l = q_l \mathbf{P}_l \begin{bmatrix} 1 \\ \tilde{\mathbf{v}}_l \end{bmatrix}, \quad (38)$$

where q_l is a non-zero scalar, \mathbf{P}_l is a $|\mathcal{T}_l| \times |\mathcal{T}_l|$ permutation matrix and $\bar{\mathbf{v}}_l$ is a vector of size $(|\mathcal{T}_l| - 1) \times 1$. Also, for any two distinct pairs $l \neq j$, the channel gain vector \mathbf{h}_{jl} can be permuted as $\bar{\mathbf{h}}_{jl} = \mathbf{P}_l^{-1} \mathbf{h}_{jl}$, and we can partition $\bar{\mathbf{h}}_{jl}$ as

$$\bar{\mathbf{h}}_{jl} = \begin{bmatrix} \bar{h}_{jl}^{(1)} \\ \bar{\mathbf{h}}_{jl}^{(2)} \end{bmatrix}, \quad (39)$$

where $\bar{h}_{jl}^{(1)}$ is a scalar and $\bar{\mathbf{h}}_{jl}^{(2)}$ is of size $(|\mathcal{T}_l| - 1) \times 1$. Therefore, the nulling condition in (36) can be rewritten as

$$\begin{bmatrix} \bar{h}_{jl}^{(1)} \bar{\mathbf{h}}_{jl}^{(2)} \\ \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \bar{\mathbf{v}}_l \end{bmatrix} = 0 \Leftrightarrow \bar{h}_{jl}^{(1)} + \bar{\mathbf{h}}_{jl}^{(2)T} \bar{\mathbf{v}}_l = 0. \quad (40)$$

Now, since the packet sent by the virtual transmitter $\widehat{\mathbf{T}}x_l$ is available in the caches of at most $|\mathcal{R}_l|$ receivers in the network, the interference of each transmitter should be nulled at least at $L - |\mathcal{R}_l| - 1$ unintended receivers. This implies that the free beamforming variables at transmitter l , i.e., $\bar{\mathbf{v}}_l$, should satisfy at least $L - |\mathcal{R}_l| - 1$ linear equations in the form of (40). This is not possible unless the number of equations is no greater than the number of variables, or

$$L - |\mathcal{R}_l| - 1 \leq |\mathcal{T}_l| - 1 \Rightarrow L \leq |\mathcal{T}_l| + |\mathcal{R}_l|. \quad (41)$$

Since the above inequality holds for all $l \in [L]$, the proof is complete. \square

B. Integer Program Formulation

Equipped with Lemma 3, we define a set of packets \mathcal{D}_m selected to be transmitted at block m to be *feasible* if its size satisfies condition (31) in Lemma 3. We can then write the following integer program (P1) to minimize the number of required communication blocks for any given caching realization and set of receiver demands:

$$\min \quad H \quad (P1-1)$$

$$\text{s.t.} \quad \bigcup_{m=1}^H \mathcal{D}_m = \bigcup_{j=1}^{K_R} (W_{d_j} \setminus \mathcal{Q}_j) \quad (P1-2)$$

$$\mathcal{D}_m \text{ is feasible, } \forall m \in [H], \quad (P1-3)$$

where (P1-2) states that all the demanded packets that are not cached at the requesting receivers need to be delivered by the transmitters over the H blocks of communication.

C. Relaxing Worst-Case Demands to Average Demands and Optimizing over Caching Realizations

We can now write an optimization problem to minimize the number of communication blocks required for delivering the worst-case demands optimized over the caching realizations. However, before that, we need to introduce some notation.

Given any caching realization $(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R})$, we can break each file $W_n, n \in [N]$, in the library into $(2^{K_T} - 1)(2^{K_R})$ subfiles $\{W_{n,\mathcal{T},\mathcal{R}}\}_{\mathcal{T} \subseteq \emptyset [K_T], \mathcal{R} \subseteq [K_R]}$, where $W_{n,\mathcal{T},\mathcal{R}}$ denotes the subfile of W_n exclusively stored in the caches of the transmitters in \mathcal{T} and receivers in \mathcal{R} , and we use the shorthand notation $\mathcal{T} \subseteq \emptyset [K_T]$ to denote $\mathcal{T} \subseteq [K_T], \mathcal{T} \neq \emptyset$. We define $a_{n,\mathcal{T},\mathcal{R}}$ as the number of packets in $W_{n,\mathcal{T},\mathcal{R}}$.

Denoting the answer to the optimization problem (P1) by $H^*(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}, \mathbf{d})$, the below optimization problem yields the number of communication blocks required for delivering the worst-case demands, minimized over all caching realizations:

$$\min_{\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}} \quad \max_{\mathbf{d}} \quad H^*(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}, \mathbf{d}) \quad (P2-1)$$

$$\text{s.t.} \quad \sum_{\mathcal{T} \subseteq \emptyset [K_T]} \sum_{\mathcal{R} \subseteq [K_R]} a_{n,\mathcal{T},\mathcal{R}} = F, \quad \forall n \in [N] \quad (P2-2)$$

$$\sum_{n=1}^N \sum_{\mathcal{R} \subseteq [K_R]} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ i \in \mathcal{T}}} a_{n,\mathcal{T},\mathcal{R}} \leq M_T F, \quad \forall i \in [K_T] \quad (P2-3)$$

$$\sum_{n=1}^N \sum_{\mathcal{T} \subseteq \emptyset [K_T]} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ j \in \mathcal{R}}} a_{n,\mathcal{T},\mathcal{R}} \leq M_R F, \quad \forall j \in [K_R] \quad (P2-4)$$

$$a_{n,\mathcal{T},\mathcal{R}} \geq 0, \quad \forall n \in [N], \forall \mathcal{T} \subseteq \emptyset [K_T], \forall \mathcal{R} \subseteq [K_R]. \quad (P2-5)$$

To lower bound the value of the above optimization problem, we can write the following optimization problem, which yields the number of communication blocks averaged over all the $\pi(N, K_R) = \frac{N!}{(N-K_R)!}$ permutations of distinct receiver demands, denoted by \mathcal{P}_{N, K_R} :

$$\min_{\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}} \frac{1}{\pi(N, K_R)} \sum_{\mathbf{d} \in \mathcal{P}_{N, K_R}} H^* \left(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}, \mathbf{d} \right) \quad (\text{P3-1})$$

$$\text{s.t.} \quad \sum_{\mathcal{T} \subseteq_{\emptyset} [K_T]} \sum_{\mathcal{R} \subseteq [K_R]} a_{n, \mathcal{T}, \mathcal{R}} = F, \quad \forall n \in [N] \quad (\text{P3-2})$$

$$\sum_{n=1}^N \sum_{\mathcal{R} \subseteq [K_R]} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ i \in \mathcal{T}}} a_{n, \mathcal{T}, \mathcal{R}} \leq M_T F, \quad \forall i \in [K_T] \quad (\text{P3-3})$$

$$\sum_{n=1}^N \sum_{\mathcal{T} \subseteq_{\emptyset} [K_T]} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ j \in \mathcal{R}}} a_{n, \mathcal{T}, \mathcal{R}} \leq M_R F, \quad \forall j \in [K_R] \quad (\text{P3-4})$$

$$a_{n, \mathcal{T}, \mathcal{R}} \geq 0, \quad \forall n \in [N], \forall \mathcal{T} \subseteq_{\emptyset} [K_T], \forall \mathcal{R} \subseteq [K_R]. \quad (\text{P3-5})$$

D. Lower Bound on the Number of Communication Blocks

Having the optimization problem in (P3), we now present the following lemma which provides a lower bound on the value of (P3).

Lemma 4. *The value of the optimization problem (P3) is bounded from below by $\frac{K_R N F (1 - \frac{M_R}{N})^2}{K_T M_T + K_R M_R}$.*

Proof. See Appendix B. □

Since the total number of packets delivered over the channel is $K_R (1 - \frac{M_R}{N}) F$ in the optimization problem (P3), Lemma 4 immediately yields the following upper bound on the one-shot linear sum-DoF:

$$\text{DoF}_{\text{L,sum}}^*(N, M_T, M_R) \leq \frac{K_R (1 - \frac{M_R}{N}) F}{\frac{K_R N F (1 - \frac{M_R}{N})^2}{K_T M_T + K_R M_R}} = \frac{K_T M_T + K_R M_R}{N - M_R}.$$

Combining the above bound with the trivial bound on the one-shot linear sum-DoF which is the number of receivers, K_R , we have

$$\text{DoF}_{\text{L,sum}}^*(N, M_T, M_R) \leq \min \left\{ \frac{K_T M_T + K_R M_R}{N - M_R}, K_R \right\}. \quad (42)$$

Now, consider the following two cases:

- $M_R \leq \frac{N}{2}$: In this case, (42) implies that

$$\begin{aligned} \text{DoF}_{\text{L,sum}}^*(N, M_T, M_R) &\leq \min \left\{ \frac{K_T M_T + K_R M_R}{N - \frac{N}{2}}, K_R \right\} \\ &\leq \min \left\{ 2 \frac{K_T M_T + K_R M_R}{N}, K_R \right\}. \end{aligned}$$

- $M_R > \frac{N}{2}$: In this case, (11) implies that one-shot linear sum-DoF of

$$\text{DoF}_{\text{L,sum}}(N, M_T, M_R) > \min \left\{ \frac{K_T M_T + K_R \frac{N}{2}}{N}, K_R \right\} > \frac{K_R}{2},$$

can be achieved, while the upper bound in (42) implies that $\text{DoF}_{\text{L,sum}}^*(N, M_T, M_R) \leq K_R$.

Therefore, in both cases, the inner bound in (11) is within a factor of 2 of the outer bound in (11), which completes the proof of the converse of Theorem 1.

VI. CONCLUDING REMARKS AND FUTURE DIRECTIONS

In this work, we considered a wireless network setting with arbitrary numbers of transmitters and receivers, where all transmitters and receivers in the network are equipped with cache memories of specific sizes. We characterized the one-shot linear sum-DoF of the network to within a gap of 2. In particular, we showed that the one-shot linear sum-DoF of the network is proportional to the aggregate cache size in the network, even though the cache of each node is isolated from all the other

nodes. We presented an achievable scheme which loads the caches carefully in order to maximize the opportunity for zero-forcing the outgoing interference from the transmitters and interference cancellation due to previously-cached content at the receivers. We also demonstrated that the achievable one-shot linear sum-DoF of our scheme is within a multiplicative factor of 2 of the optimal one-shot linear sum-DoF by bounding the number of communication blocks required to deliver any set of requested files to the receivers using an integer programming approach.

There are several interesting directions following this work. First, in this work we assumed all the links in the network to be present in the network topology. However, due to fading effects, some links between certain transmitter-receiver pairs might be absent from the network topology. It would be interesting to study what type of caching strategies are optimal in this case and to explore its connections to the index coding problem [16–18]. Another direction would be to combine caching with more sophisticated interference management schemes. Some initial results have been reported in [19], in which the authors used the replication in the cache contents at the transmitters in order to improve the system performance using the ITLinQ scheme [20–22]. It would be interesting to study the role of transmitter and receiver caches illustrated in this work in improving the achievable system throughput that more sophisticated delivery schemes such as ITLinQ can provide.

APPENDIX A PROOF OF LEMMA 1

For any $\mathcal{T} \subseteq [K_T]$ s.t. $|\mathcal{T}| = t_T$ and for any $l \in [t_T + t_R]$, it is clear that the set $\mathcal{T} \oplus_{K_T} (l-1)$ is of size t_T . Also, for any $\mathcal{R} \subseteq [K_R]$ s.t. $|\mathcal{R}| = t_T + t_R$, and for any permutation $\pi \in \Pi_{\mathcal{R}}^{\text{circ}}$, the vector $\pi[l+1 : l+t_R]$ is of size t_R and the vector $\pi[l+t_R+1 : l+t_R+t_T-1]$ is of size t_T-1 .

Furthermore, note that $W_{d_{\pi(l)}, \mathcal{T} \oplus_{K_T} (l-1), \pi[l+1:l+t_R], \pi[l+t_R+1:l+t_R+t_T-1]}$ is a subfile of the file $W_{d_{\pi(l)}}$ requested by receiver $\text{Rx}_{\pi(l)}$. However, since $\pi(l) \notin \pi[l+1 : l+t_R]$, receiver $\text{Rx}_{\pi(l)}$ has not stored the packets in this subfile in its cache and therefore, this subfile needs to be delivered to this receiver.

Finally, each set inside the union in (19) is composed of $t_T + t_R$ subfiles. The number of such sets is equal to

$$\binom{K_T}{t_T} \binom{K_R}{t_T + t_R} (t_T + t_R - 1)! \quad (43)$$

Hence, the total number of subfiles in (19) is equal to

$$\binom{K_T}{t_T} \binom{K_R}{t_T + t_R} (t_T + t_R - 1)! (t_T + t_R) = \binom{K_T}{t_T} \binom{K_R}{t_T + t_R} (t_T + t_R)! \quad (44)$$

On the other hand, each receiver Rx_j has already cached $\binom{K_T}{t_T} \binom{K_R-1}{t_R-1}$ subfiles as in (16) in its cache, and needs the rest of the subfiles of its requested file, i.e., $\binom{K_T}{t_T} \binom{K_R-1}{t_R}$ subfiles, where each subfile is further partitioned into $\frac{t_R! [K_R - (t_R + 1)]!}{[K_R - (t_R + t_T)]!}$ smaller subfiles. Hence, the total number of small subfiles that need to be delivered to all the receivers is equal to

$$K_R \left[\binom{K_T}{t_T} \binom{K_R-1}{t_R} \right] \left[\frac{t_R! [K_R - (t_R + 1)]!}{[K_R - (t_R + t_T)]!} \right] = \binom{K_T}{t_T} \binom{K_R}{t_T + t_R} (t_T + t_R)!, \quad (45)$$

which equals the total number of small subfiles in (19), calculated in (44). Consequently, the set of requested subfiles which are not cached at the corresponding receivers can be partitioned as in (19), hence the proof is complete. \square

APPENDIX B PROOF OF LEMMA 4

According to the constraint (31), each of the packets of order s , which are available at s nodes, either on the transmitter side or the receiver side, can be scheduled with at most $s-1$ packets of the same order. Therefore, for any given caching realization and set of demands, we have the lower bound

$$\begin{aligned} & H^* \left(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}, \mathbf{d} \right) \\ & \geq \sum_{s=K_R}^{K_T+K_R} \sum_{j=1}^{K_R} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}| \in [s]}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}| = s - |\mathcal{S}_T| \\ j \notin \mathcal{R}}} \frac{a_{d_j, \mathcal{T}, \mathcal{R}}}{K_R} + \sum_{s=1}^{K_R-1} \sum_{j=1}^{K_R} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}| \in [s]}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}| = s - |\mathcal{S}_T| \\ j \notin \mathcal{R}}} \frac{a_{d_j, \mathcal{T}, \mathcal{R}}}{s} \\ & \geq \sum_{s=1}^{K_T+K_R} \sum_{j=1}^{K_R} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}| \in [s]}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}| = s - |\mathcal{S}_T| \\ j \notin \mathcal{R}}} \frac{a_{d_j, \mathcal{T}, \mathcal{R}}}{s}. \end{aligned} \quad (46)$$

Now, denoting the objective function in (P3-1) by $\bar{H}\left(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}\right)$, we have

$$\begin{aligned}
\bar{H}\left(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R}\right) &\geq \frac{1}{\pi(N, K_R)} \sum_{s=1}^{K_T+K_R} \frac{1}{s} \sum_{j=1}^{K_R} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}| \in [s]}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=s-|S_{\mathcal{T}}| \\ j \notin \mathcal{R}}} \pi(N-1, K_R-1) \sum_{n=1}^N a_{n, \mathcal{T}, \mathcal{R}} \\
&= \frac{1}{N} \sum_{s=1}^{K_T+K_R} \frac{1}{s} \sum_{j=1}^{K_R} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}| \in [s]}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=s-|S_{\mathcal{T}}| \\ j \notin \mathcal{R}}} \sum_{n=1}^N a_{n, \mathcal{T}, \mathcal{R}} \\
&= \frac{1}{N} \sum_{r=1}^{K_T} \sum_{r'=0}^{K_R} \frac{1}{r+r'} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r' \\ j \notin \mathcal{R}}} \sum_{n=1}^N a_{n, \mathcal{T}, \mathcal{R}} \\
&= \frac{1}{N} \sum_{r=1}^{K_T} \sum_{r'=0}^{K_R} \frac{K_R-r'}{r+r'} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r'}} \sum_{n=1}^N a_{n, \mathcal{T}, \mathcal{R}} \\
&= \frac{1}{N} \sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} \frac{b_{r, r'}}{r+r'}, \tag{47}
\end{aligned}$$

where for any $r \in [K_T]$ and $r' \in [K_R-1] \cup \{0\}$, we define

$$b_{r, r'} \triangleq \sum_{j=1}^{K_R} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r' \\ j \notin \mathcal{R}}} \sum_{n=1}^N a_{n, \mathcal{T}, \mathcal{R}} = (K_R-r') \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r'}} \sum_{n=1}^N a_{n, \mathcal{T}, \mathcal{R}}. \tag{48}$$

Moreover, adding the constraint in (P3-3) over all transmitters yields

$$K_T M_T F \geq \sum_{i=1}^{K_T} \sum_{n=1}^N \sum_{\mathcal{R} \subseteq [K_R]} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ i \in \mathcal{T}}} a_{n, \mathcal{T}, \mathcal{R}} \tag{49}$$

$$= \sum_{n=1}^N \sum_{\mathcal{R} \subseteq [K_R]} \sum_{i=1}^{K_T} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ i \in \mathcal{T}}} a_{n, \mathcal{T}, \mathcal{R}} \tag{50}$$

$$= \sum_{n=1}^N \sum_{\mathcal{R} \subseteq [K_R]} \sum_{r=1}^{K_T} r \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} a_{n, \mathcal{T}, \mathcal{R}}. \tag{51}$$

Likewise, adding the constraint in (P3-4) over all receivers yields

$$K_R M_R F \geq \sum_{j=1}^{K_R} \sum_{n=1}^N \sum_{\mathcal{T} \subseteq \emptyset [K_T]} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ j \in \mathcal{R}}} a_{n, \mathcal{T}, \mathcal{R}} \tag{52}$$

$$= \sum_{n=1}^N \sum_{\mathcal{T} \subseteq \emptyset [K_T]} \sum_{j=1}^{K_R} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ j \in \mathcal{R}}} a_{n, \mathcal{T}, \mathcal{R}} \tag{53}$$

$$= \sum_{n=1}^N \sum_{\mathcal{T} \subseteq \emptyset [K_T]} \sum_{r'=0}^{K_R} r' \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r'}} a_{n, \mathcal{T}, \mathcal{R}}, \tag{54}$$

and from (51) and (54), we have

$$(K_T M_T + K_R M_R)F \geq \sum_{n=1}^N \left[\sum_{\mathcal{R} \subseteq [K_R]} \sum_{r=1}^{K_T} r \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} a_{n,\mathcal{T},\mathcal{R}} + \sum_{\mathcal{T} \subseteq [K_T]} \sum_{r'=0}^{K_R} r' \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r'}} a_{n,\mathcal{T},\mathcal{R}} \right] \quad (55)$$

$$= \sum_{r=1}^{K_T} \sum_{r'=0}^{K_R} (r+r') \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r'}} \sum_{n=1}^N a_{n,\mathcal{T},\mathcal{R}} \quad (56)$$

$$\geq \sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} \frac{r+r'}{K_R-r'} b_{r,r'}. \quad (57)$$

Now, using the Cauchy-Schwarz inequality, we can write

$$\sum_{r'=0}^{K_R-1} b_{r,r'} \leq \sqrt{\sum_{r'=0}^{K_R-1} \frac{r+r'}{K_R-r'} b_{r,r'}} \sqrt{\sum_{r'=0}^{K_R-1} \frac{K_R-r'}{r+r'} b_{r,r'}}. \quad (58)$$

Summing the above inequality over r yields

$$\sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} b_{r,r'} \leq \sum_{r=1}^{K_T} \left[\sqrt{\sum_{r'=0}^{K_R-1} \frac{r+r'}{K_R-r'} b_{r,r'}} \sqrt{\sum_{r'=0}^{K_R-1} \frac{K_R-r'}{r+r'} b_{r,r'}} \right] \quad (59)$$

$$\leq \sqrt{\sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} \frac{r+r'}{K_R-r'} b_{r,r'}} \sqrt{\sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} \frac{K_R-r'}{r+r'} b_{r,r'}} \quad (60)$$

$$\leq \sqrt{(K_T M_T + K_R M_R)F} \sqrt{\sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} \frac{K_R-r'}{r+r'} b_{r,r'}}, \quad (61)$$

where in (60) we have invoked the Cauchy-Schwarz inequality again and (61) follows from (57). On the other hand, we have

$$\sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} b_{r,r'} = \sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} \sum_{j=1}^{K_R} \sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r' \\ j \notin \mathcal{R}}} \sum_{n=1}^N a_{n,\mathcal{T},\mathcal{R}} \quad (62)$$

$$= \sum_{r=1}^{K_T} \sum_{r'=0}^{K_R} \sum_{j=1}^{K_R} \left[\left(\sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r'}} \sum_{n=1}^N a_{n,\mathcal{T},\mathcal{R}} \right) - \left(\sum_{\substack{\mathcal{T} \subseteq [K_T]: \\ |\mathcal{T}|=r}} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ |\mathcal{R}|=r' \\ j \in \mathcal{R}}} \sum_{n=1}^N a_{n,\mathcal{T},\mathcal{R}} \right) \right] \quad (63)$$

$$= K_R \left(\sum_{n=1}^N \sum_{\mathcal{T} \subseteq \emptyset [K_T]} \sum_{\mathcal{R} \subseteq [K_R]} a_{n,\mathcal{T},\mathcal{R}} \right) - \sum_{j=1}^{K_R} \sum_{n=1}^N \sum_{\mathcal{T} \subseteq \emptyset [K_T]} \sum_{\substack{\mathcal{R} \subseteq [K_R]: \\ j \in \mathcal{R}}} a_{n,\mathcal{T},\mathcal{R}} \quad (64)$$

$$\geq K_R(N - M_R)F, \quad (65)$$

where the inequality is due to (P3-2) and (52). Therefore, we can continue (47) to bound the objective function in (P3-1) as

$$\bar{H} \left(\{\mathcal{P}_i\}_{i=1}^{K_T}, \{\mathcal{Q}_i\}_{i=1}^{K_R} \right) \geq \frac{1}{N} \sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} \frac{b_{r,r'}}{r+r'} \quad (66)$$

$$\geq \frac{1}{K_R N} \sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} \frac{(K_R-r')b_{r,r'}}{r+r'} \quad (67)$$

$$\geq \frac{1}{K_R N F (K_T M_T + K_R M_R)} \left(\sum_{r=1}^{K_T} \sum_{r'=0}^{K_R-1} b_{r,r'} \right)^2 \quad (68)$$

$$\geq \frac{1}{K_R N F (K_T M_T + K_R M_R)} \left(K_R (N - M_R) F \right)^2 \quad (69)$$

$$= \frac{K_R N F \left(1 - \frac{M_R}{N}\right)^2}{K_T M_T + K_R M_R}, \quad (70)$$

where (68) and (69) follow from (61) and (65), respectively. This completes the proof. \square

REFERENCES

- [1] N. Golrezaei, K. Shanmugam, A. G. Dimakis, A. F. Molisch, and G. Caire, "Femtocaching: Wireless video content delivery through distributed caching helpers," in *INFOCOM*, March 2012, pp. 1107–1115.
- [2] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *Information Theory, IEEE Transactions on*, vol. 60, no. 5, pp. 2856–2867, 2014.
- [3] S. P. Shariatpanahi, S. A. Motahari, and B. H. Khalaj, "Multi-server coded caching," *arXiv preprint arXiv:1503.00265*, 2015.
- [4] A. Liu and V. K. N. Lau, "Exploiting base station caching in MIMO cellular networks: Opportunistic cooperation for video streaming," *Signal Processing, IEEE Transactions on*, vol. 63, no. 1, pp. 57–69, Jan. 2015.
- [5] A. Sengupta, R. Tandon, and O. Simeone, "Cache aided wireless networks: Tradeoffs between storage and latency," *arXiv preprint arXiv:1512.07856*, 2015.
- [6] M. Ji, G. Caire, and A. F. Molisch, "Wireless device-to-device caching networks: Basic principles and system performance," *arXiv preprint arXiv:1305.5216*, 2013.
- [7] M. A. Maddah-Ali and U. Niesen, "Cache-aided interference channels," in *Information Theory (ISIT), 2015 IEEE International Symposium on*. IEEE, 2015, pp. 809–813.
- [8] B. Azari, O. Simeone, U. Spagnolini, and A. Tulino, "Hypergraph-based analysis of clustered cooperative beamforming with application to edge caching," *arXiv preprint arXiv:1510.06222*, 2015.
- [9] S.-H. Park, O. Simeone, and S. Shamai, "Joint optimization of cloud and edge processing for fog radio access networks," *arXiv preprint arXiv:1601.02460*, 2016.
- [10] M. Afshang, H. S. Dhillon, and P. H. J. Chong, "Fundamentals of cluster-centric content placement in cache-enabled device-to-device networks," *arXiv preprint arXiv:1509.04747*, 2015.
- [11] M. Tao, E. Chen, H. Zhou, and W. Yu, "Content-centric sparse multicast beamforming for cache-enabled cloud RAN," *arXiv preprint arXiv:1512.06938*, 2015.
- [12] Y. Ugur, Z. H. Awan, and A. Sezgin, "Cloud radio access networks with coded caching," *arXiv preprint arXiv:1512.02385*, 2015.
- [13] X. Peng, J.-C. Shen, J. Zhang, and K. B. Letaief, "Backhaul-aware caching placement for wireless networks," *arXiv preprint arXiv:1509.00558*, 2015.
- [14] M. Razaviyayn, G. Lyubeznik, and Z.-Q. Luo, "On the degrees of freedom achievable through interference alignment in a MIMO interference channel," *Signal Processing, IEEE Transactions on*, vol. 60, no. 2, pp. 812–821, 2012.
- [15] H. Weingarten, Y. Steinberg, and S. Shamai, "The capacity region of the gaussian multiple-input multiple-output broadcast channel," *Information Theory, IEEE Transactions on*, vol. 52, no. 9, pp. 3936–3964, Sept 2006.
- [16] Y. Birk and T. Kol, "Informed-source coding-on-demand (ISCOD) over broadcast channels," in *INFOCOM '98. Seventeenth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, vol. 3, Mar 1998, pp. 1257–1264.
- [17] Z. Bar-Yossef, Y. Birk, T. Jayram, and T. Kol, "Index coding with side information," in *Foundations of Computer Science, 2006. FOCS '06. 47th Annual IEEE Symposium on*, Oct 2006, pp. 197–206.
- [18] S. El Rouayheb, A. Sprintson, and C. Georghiadis, "On the index coding problem and its relation to network coding and matroid theory," *Information Theory, IEEE Transactions on*, vol. 56, no. 7, pp. 3187–3195, 2010.
- [19] N. Naderializadeh, D. T. Kao, and A. S. Avestimehr, "How to utilize caching to improve spectral efficiency in device-to-device wireless networks," in *Communication, Control, and Computing (Allerton), 2014 52nd Annual Allerton Conference on*. IEEE, 2014, pp. 415–422.
- [20] N. Naderializadeh and A. S. Avestimehr, "ITLinQ: A new approach for spectrum sharing in device-to-device communication systems," *Selected Areas in Communications, IEEE Journal on*, vol. 32, no. 6, pp. 1139–1151, June 2014.
- [21] —, "ITLinQ: A new approach for spectrum sharing," in *Dynamic Spectrum Access Networks (DYSPAN), 2014 IEEE International Symposium on*. IEEE, 2014, pp. 327–333.
- [22] —, "ITLinQ: A new approach for spectrum sharing in device-to-device communication systems," in *Information Theory (ISIT), 2014 IEEE International Symposium on*, June 2014, pp. 1573–1577.