

# The Optimal Memory-Rate Trade-off for the Non-uniform Centralized Caching Problem with Two Files under Uncoded Placement

Saeid Sahraei, *Member, IEEE*, Pierre Quinton, *Student Member, IEEE*, and Michael Gastpar, *Fellow, IEEE*

**Abstract**—A new scheme for the problem of centralized coded caching with non-uniform demands is proposed. The distinguishing feature of the proposed placement strategy is that it admits equal sub-packetization for all files while allowing the users to allocate more cache to the files which are more popular. This creates natural broadcasting opportunities in the delivery phase which are simultaneously helpful for the users who have requested files of different popularities. For the case of two files, we propose a new delivery strategy based on interference alignment which enables each user to decode his desired file following a two-layer peeling decoder. Furthermore, we extend the existing converse bounds for uniform demands under uncoded placement to the nonuniform case. To accomplish this, we construct  $N!$  auxiliary users, corresponding to all permutations of the  $N$  files, each caching carefully selected sub-packets of the files. Each auxiliary user provides a different converse bound. The overall converse bound is the maximum of all these  $N!$  bounds. We prove that our achievable delivery rate for the case of two files meets this converse, thereby establishing the optimal expected memory-rate trade-off for the case of  $K$  users and two files with arbitrary popularities under uncoded placement.

**Index Terms**—Coded Caching, Non-Uniform Demands, Combinatorial Design

## I. INTRODUCTION

WIRELESS traffic has been dramatically increasing in recent years, mainly due to the increasing popularity of video streaming services. Caching is a mechanism for Content Distribution Networks (CDNs) to cope with this increasing demand by placing the contents closer to the users during off-peak hours. The attractive possibility of replacing expensive bandwidth with cheap memories has caused an outburst of research in the recent years [1]–[9]. Coded caching [1] is a canonical formulation of a two-stage communication problem between a server and many clients which are connected to the server via a shared broadcast channel. The two stages consist of filling in the caches of the users during off-peak hours and transmitting the desired data to them at their request, typically

during the peak hours. The local caches of the users act as side information for an instance of the index coding problem where different users may have different demands. Logically, if a content is more likely to be requested, it is desirable to cache more of it during the first stage. Furthermore, by diversifying the contents cached at different users, broadcasting opportunities can be created which are simultaneously helpful for several users [1].

In general, there exists a trade-off between the amount of cache that each user has access to and the delivery rate at the second stage. Significant progress has been made towards characterizing this trade-off for worst case and average case demands under uniform file popularities [1]–[3], [10], [11]. The optimal memory-rate region under uncoded placement is known [3], [12], and the general optimal memory-rate region has been characterized within a factor of 2 [2]. Furthermore, many achievability results have been proposed based on coded placement [11], [13], [14]. Some of these schemes outperform the optimal caching scheme under uncoded placement [3], establishing that uncoded placement is in general sub-optimal.

By contrast, the coded caching problem with non-uniform file popularities, an arguably more realistic model, has remained largely open. The existing achievability schemes are generally speaking straightforward generalizations of the caching schemes that are specifically tailored to the uniform case. Here we briefly review some of these works.

### A. Related Work

The main body of work on non-uniform coded caching has been concentrated around the decentralized paradigm where there is no coordination among different users [10]. The core idea here is to partition the files into  $L$  groups where the files within each group have similar popularities [15]. Within each group, one performs the decentralized coded caching strategy of [10] as if all the files had the same probability of being requested. In the delivery phase, coding opportunities among different groups are ignored and as a result, the total delivery rate is the sum of the delivery rates for the  $L$  partitions. It was subsequently suggested to use  $L = 2$  groups [16]–[18] and to allocate no cache at all to the group which contains the least popular files. This simple scheme was proven to be within a multiplicative and additive gap of optimal for arbitrary file popularities [18].

The problem of *centralized* coded caching with non-uniform demands has also been extensively studied [19]–[23]. Here, in

This work was supported in part by the Swiss National Science Foundation under Grants 169294 and 178309. This paper was partially presented at the International Zurich Seminar on Information and Communication, 2018.

Saeid Sahraei is with the Department of Electrical Engineering, University of Southern California, Los Angeles, CA 90089, USA (e-mail: ss\_805@usc.edu).

Pierre Quinton and Michael Gastpar are with the School of Computer and Communication Sciences, École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Vaud 1015, Switzerland (e-mail: pierre.quinton@epfl.ch, michael.gastpar@epfl.ch).

Copyright (c) 2019 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

order to create coding opportunities among files with varying popularities, a different approach has been taken. Each file is partitioned into  $2^K$  subfiles corresponding to all possible ways that a given subfile can be shared among  $K$  users. This creates coding opportunities among subfiles that are shared among equal number of users, even if they belong to files with different popularities. The delivery rate can be minimized by solving an optimization problem that decides what portion of each file must be shared among  $i$  users, for any  $i \in [0 : K]$ . It was proven in [21] that if the size of the cache belongs to a set of base-cases  $\mathcal{M}$  of size  $NK + 1$ , the best approach is to allocate no cache at all to the least popular files while treating the other files as if they were equally probable of being requested. Memory-sharing among such points must be performed if the cache size is not a member of  $\mathcal{M}$ . The set of base-cases depends on the popularities of the files, the number of users and the number of files, and can be computed via an efficient algorithm [21].

The graph-based delivery strategies that are predominantly used in this line of research [17], [21] are inherently limited, in that they do not capture the algebraic properties of summation over finite fields. For instance, one can easily construct examples where the chromatic-number index coding scheme in [17] is sub-optimal, even for uniform file popularities. Suppose we have only one file  $A = \{A_1, A_2, A_3\}$  and 3 users each with a cache of size  $1/3$  file. Assume user  $i$  caches  $A_i$ . In this case, the optimal delivery rate is  $2/3$  but the clique-cover (or chromatic-number) approach in [17] only provides a delivery rate of 1. This is due to the fact that from  $A_1 \oplus A_2$  and  $A_2 \oplus A_3$  one can recover  $A_1 \oplus A_3$ , a property that the graph-based model in [17] fails to reflect. This issue was addressed in a systematic manner by Yu et. al. in [3] which introduced the concept of leaders. Our delivery scheme in this paper provides an alternative mechanism for overcoming this limitation, which transcends the uniform file popularities and can be applied to nonuniform demands. This is accomplished via interference alignment as outlined in the proof of achievability of Lemma 1 in Section VI.

In [23] it was proven that a slightly modified version of the decentralized scheme in [16] is optimal for the centralized caching problem when we have only two users. In [22], a centralized caching strategy was proposed for the case where the number of users  $K$  is prime and the case where  $K$  divides  $\binom{K}{t}$ , where  $\binom{K}{t}$  is the subpacketization of each file. The placement scheme allows for equal subpacketization of all the files while more stringent requirements are imposed for caching subfiles of less popular files. This concept is closely related to what was presented in [24] which serves as the placement scheme for the current paper.

## B. Our Contributions

In this paper, we first propose a centralized placement strategy for an arbitrary number of users and files, which allows for equal subpacketization of all the files while allocating less cache to the files which are less likely to be requested. This creates natural coding opportunities in the delivery phase among all the files regardless of their popularities. Next, we

propose a delivery strategy for the case of two files and an arbitrary number of users. This delivery strategy consists of two phases. First, each file is compressed down to its entropy conditioned on the side information available at the users who have requested it. Simultaneously, this encoding aims at *aligning* the subfiles which are unknown to the users who have *not* requested them. In the second phase of the delivery strategy, the two encoded files are further encoded with an MDS code and broadcast to the users. Each user will be able to decode his desired file following a two-layer peeling decoder. By extending the converse bound for uncoded placement first proposed in [3] to the non-uniform case, we prove that our joint placement and delivery strategy is optimal for two files with arbitrary popularities under uncoded placement. To summarize, our main contributions are the following:

- A new placement strategy is developed for non-uniform caching with  $K$  users and  $N$  files (Section V). This scheme allows for equal sub-packetization of every file, while allocating more cache to files that are more popular. A simple modification of the proposed scheme can be applied to user-dependent file popularities. More broadly, the proposed multiset indexing approach to subpacketization can be expected to find applications in other coding problems of combinatorial nature with heterogeneous objects, such as Coded Data shuffling [25], Coded Map-Reduce [26], and Fractional Repetition codes [27] for Distributed Storage.
- An extension of the converse bound under uncoded placement first proposed in [3] to non-uniform caching with  $K$  users and  $N$  files is established (Section VII).
- A new delivery strategy is presented for the case of two files which relies on source coding and interference alignment (Section VI). The achievable expected delivery rate meets the extended converse bound under uncoded placement, hence establishing the optimal memory-rate tradeoff for non-uniform demands for the case of two files. If each file has probability  $1/2$ , this approach leads to an alternative delivery strategy for uniform caching of two files, which can be of independent interest.

The rest of the paper is organized as follows. We introduce the notation used throughout the paper and the formal problem statement in Sections II and III. In Section IV we will explain the main ideas behind our caching strategy via a case study. The general placement and delivery strategy are presented in Sections V and VI. We will then propose our converse bound under uncoded placement in Section VII. In Section VIII we will prove that our proposed caching strategy is optimal for the case of two files. Finally, in Section IX, we will provide numerical results, and conclude the paper in Section X.

## II. NOTATION

For two integers  $a, b$  define  $\binom{a}{b} = 0$  if  $b < 0$  or  $b > a$ . For  $n+1$  non-negative integers  $a, b_1, \dots, b_n$  that satisfy  $\sum_{i=1}^n b_i = a$ , define

$$\binom{a}{b_1, \dots, b_n} = \binom{a}{b_1} \binom{a-b_1}{b_2} \dots \binom{a-\sum_{i=1}^{n-1} b_i}{b_n}$$

$$= \frac{a!}{b_1! \cdots b_n!}. \quad (1)$$

For a positive integer  $a$  define  $[a] = \{1, \dots, a\}$ . For two integers  $a, b$  define  $[a : b] = \{a, a+1, \dots, b\}$ . For two column vectors  $\mathbf{u}$  and  $\mathbf{v}$  denote their vertical concatenation by  $[\mathbf{u}; \mathbf{v}]$ . For a real number  $a$ , define  $\lfloor a \rfloor$  as the largest integer no greater than  $a$ . Similarly, define  $\lceil a \rceil$  as the smallest integer no less than  $a$ . For  $q \in \mathbb{R}^+$  and a discrete random variable  $X$  with support  $\mathcal{X}$  define  $H_q(X)$  as the entropy of  $X$  in base  $q$ :

$$H_q(X) = - \sum_{x \in \mathcal{X}} \mathbb{P}(X = x) \log_q \mathbb{P}(X = x). \quad (2)$$

Suppose we have a function  $f(\cdot) : \mathcal{D} \rightarrow \mathbb{R}$  where  $\mathcal{D}$  is a discrete set of points in  $\mathbb{R}^n$ . Let  $\mathcal{T}$  be the convex hull of  $\mathcal{D}$ . Define

$$\begin{aligned} g(\mathbf{t}) &= \mathcal{L}_{\mathbf{r} \rightarrow \mathbf{t}} f(\mathbf{r}) \\ g(\cdot) &: \mathcal{T} \rightarrow \mathbb{R} \end{aligned} \quad (3)$$

as the lower convex envelope of  $f(\cdot)$  evaluated at point  $\mathbf{t} \in \mathcal{T}$ .

### III. MODEL DESCRIPTION AND PROBLEM STATEMENT

We follow the canonical model in [1] except here we concentrate on the expected delivery rate as opposed to the worst case delivery rate. For the sake of completeness, we repeat the model description here. We have a network consisting of  $K$  users that are connected to a server through a shared broadcast link. The server has access to  $N$  files  $W_1, \dots, W_N$  each of size  $F$  symbols over a sufficiently large field  $\mathbb{F}_q$ . Therefore,  $H_q(W_i) \leq F$ . Each user has a cache of size  $M$  symbols over  $\mathbb{F}_q$ . An illustration of the network has been provided in Figure 1. The communication between the server and the users takes place in two phases, placement and delivery.

In the placement phase, each user stores some function of all the files  $Z_i = f_i(W_1, \dots, W_N), i \in [K]$  in his local cache. Therefore, for a fixed (normalized) memory size  $M$ , a placement strategy  $\mathcal{M}$  consists of  $K$  placement functions  $Z_i = f_i(W_1, \dots, W_N), i \in [K]$  such that  $H_q(Z_i) \leq MF$  for all  $i \in [K]$ . After the placement phase, each user requests one file from the server. We represent the request of the  $i$ 'th user with  $d_i \in [N]$  which is drawn from a known probability distribution  $\mathbf{p}$ . Furthermore, the requests of all the users are independent and identically distributed. After receiving the request vector  $\mathbf{d}$ , the server transmits a delivery message  $X_{\mathbf{d}, \mathcal{M}, F}$  through the broadcast link to all the users. User  $i$  then computes a function  $\hat{W}_{d_i} = g_i(X_{\mathbf{d}, \mathcal{M}, F}, Z_i, \mathbf{d})$  in order to estimate  $W_{d_i}$ . For a fixed placement strategy  $\mathcal{M}$ , fixed file size  $F$ , and fixed request vector  $\mathbf{d}$  we say that a delivery rate of  $R_{\mathbf{d}, \mathcal{M}, F}$  is achievable if a delivery message  $X_{\mathbf{d}, \mathcal{M}, F}$  and decoding functions  $g_i(\cdot), i \in [K]$  exist such that

$$\mathbb{P}(g_i(X_{\mathbf{d}, \mathcal{M}, F}, Z_i, \mathbf{d}) \neq W_{d_i}) = 0, \forall i \in [N], \quad (4)$$

and

$$H_q(X_{\mathbf{d}, \mathcal{M}, F}) \leq R_{\mathbf{d}, \mathcal{M}, F} F. \quad (5)$$

For a fixed placement strategy  $\mathcal{M}$ , we say that an expected delivery rate  $\bar{R}_{\mathcal{M}}$  is achievable if there exists a sequence of

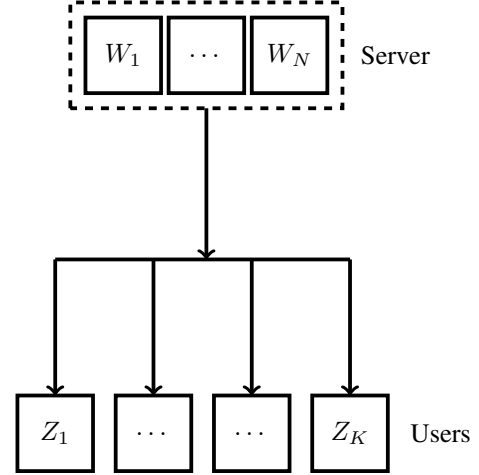


Fig. 1: An illustration of the caching network. A server is connected to  $K$  users via a shared broadcast link. Each user has a cache of size  $MF$  symbols where he can store an arbitrary function of the files  $W_1, \dots, W_N$ .

achievable delivery rates  $\{R_{\mathbf{d}, \mathcal{M}, F} | \mathbf{d} \in [N]^K, F \in \mathbb{N}\}$  such that

$$\limsup_{F \rightarrow \infty} \mathbb{E}_{\mathbf{d}} R_{\mathbf{d}, \mathcal{M}, F} \leq \bar{R}_{\mathcal{M}}. \quad (6)$$

Finally, for a memory of size  $M$ , we say that an expected delivery rate  $\bar{R}$  is achievable if there exists a placement strategy  $\mathcal{M} = (Z_1, \dots, Z_K)$  with  $H_q(Z_i) \leq MF$  for all  $i \in [K]$ , for which an expected delivery rate of  $\bar{R}_{\mathcal{M}} \leq \bar{R}$  is achievable.

Our goal in this paper is to characterize the minimum expected delivery rate for all  $M$  under the restriction of uncoded placement. In other words, the placement functions must be of the form

$$\begin{aligned} Z_i &= f_i(W_1, \dots, W_N) \\ &= (W_1|_{A_1}, \dots, W_N|_{A_N}) \text{ for all } i \in [K], \end{aligned} \quad (7)$$

where  $A_j \subseteq [F]$  for all  $j \in [N]$ , and  $W_j|_{A_j}$  refers to the subset of symbols of the file  $W_j$  which are indexed in the set  $A_j$ .

### IV. MOTIVATING EXAMPLE: THE CASE OF FOUR USERS AND TWO FILES

Consider the caching problem with two files  $W_1$  and  $W_2$  and  $K = 4$  users. Assume the probability of requesting  $W_2$  is lower than the probability of requesting  $W_1$ . In this section we will demonstrate how to find the optimal expected delivery rate for any memory size for this particular choice of parameters, while explaining the main principles behind our joint placement and delivery strategy. We start by fixing two integers  $r_1, r_2$  such that  $0 \leq r_2 \leq r_1 \leq K$ . As we will see soon, any choice of  $(r_1, r_2)$  corresponds to a particular  $(M_1, M_2)$  where  $M_i$  is the amount of cache that each user allocates to file  $W_i$ , normalized by the size of one file. For the sake of brevity, we will explain our strategy only for  $(r_1, r_2) = (2, 1)$ . The delivery rate for other possible choices

$Z_1$	$Z_2$
$W_{1,\{1,2\},\{1\}}, W_{1,\{1,2\},\{2\}}$	$W_{1,\{1,2\},\{2\}}, W_{1,\{1,2\},\{1\}}$
$W_{1,\{1,3\},\{1\}}, W_{1,\{1,3\},\{3\}}$	$W_{1,\{2,3\},\{2\}}, W_{1,\{2,3\},\{3\}}$
$W_{1,\{1,4\},\{1\}}, W_{1,\{1,4\},\{4\}}$	$W_{1,\{2,4\},\{2\}}, W_{1,\{2,4\},\{4\}}$
$W_{2,\{1,2\},\{1\}}$	$W_{2,\{1,2\},\{2\}}$
$W_{2,\{1,3\},\{1\}}$	$W_{2,\{2,3\},\{2\}}$
$W_{2,\{1,4\},\{1\}}$	$W_{2,\{2,4\},\{2\}}$
$Z_3$	$Z_4$
$W_{1,\{1,3\},\{3\}}, W_{1,\{1,3\},\{1\}}$	$W_{1,\{1,4\},\{4\}}, W_{1,\{1,4\},\{1\}}$
$W_{1,\{2,3\},\{3\}}, W_{1,\{2,3\},\{2\}}$	$W_{1,\{2,4\},\{4\}}, W_{1,\{2,4\},\{2\}}$
$W_{1,\{3,4\},\{3\}}, W_{1,\{3,4\},\{4\}}$	$W_{1,\{3,4\},\{4\}}, W_{1,\{3,4\},\{3\}}$
$W_{2,\{1,3\},\{3\}}$	$W_{2,\{1,4\},\{4\}}$
$W_{2,\{2,3\},\{3\}}$	$W_{2,\{2,4\},\{4\}}$
$W_{2,\{3,4\},\{3\}}$	$W_{2,\{3,4\},\{4\}}$

TABLE I: The proposed placement scheme for  $N = 2$ ,  $K = 4$ ,  $(r_1, r_2) = (2, 1)$ .

of  $(r_1, r_2)$  will be summarized at the end of this section. Next, we will characterize the entire  $(M_1, M_2, R)$  region that can be achieved by our algorithm. Finally, we will illustrate how to find the optimal choice of  $(M_1, M_2)$  for a particular cache size  $M$ .

Define the parameter  $S = \binom{K}{r_1} \binom{r_1}{r_2}$ . We divide each of the two files into  $S$  subfiles and index them as  $W_{1,\tau_1,\tau_2}$  and  $W_{2,\tau_1,\tau_2}$  such that  $\tau_2 \subseteq \tau_1 \subseteq [K]$  and  $|\tau_1| = r_1, |\tau_2| = r_2$ . In this example,  $S = \binom{4}{2} \binom{2}{1} = 12$ . The 12 subfiles of  $W_i$  are then denoted by  $W_{i,\{1,2\},\{1\}}, W_{i,\{1,2\},\{2\}}, W_{i,\{1,3\},\{1\}}, W_{i,\{1,3\},\{3\}}, W_{i,\{1,4\},\{1\}}, W_{i,\{1,4\},\{4\}}, W_{i,\{2,3\},\{2\}}, W_{i,\{2,3\},\{3\}}, W_{i,\{2,4\},\{2\}}, W_{i,\{2,4\},\{4\}}, W_{i,\{3,4\},\{3\}}, W_{i,\{3,4\},\{4\}}$ . In our placement strategy, user  $j$  stores the subfiles  $W_{1,\tau_1,\tau_2}$  for which  $j \in \tau_1$ , as well as the subfiles  $W_{2,\tau_1,\tau_2}$  for which  $j \in \tau_2$ . Since  $\tau_2 \subseteq \tau_1$ , the users naturally store fewer subfiles of  $W_2$  than  $W_1$ . In our running example, each user stores six subfiles of  $W_1$  but only three subfiles of  $W_2$ . The cache contents of each user has been summarized in Table I.

Note that this placement scheme results in a memory of size  $M = \frac{3}{4}$ . As we will see soon, the memory size is in general  $M = \sum_{i=1}^N M_i$  where  $M_i = \frac{r_i}{K}$  is the amount of cache dedicated by each user to file  $W_i$ . It is important to note that despite the fact that each user has allocated more cache to file  $W_1$ , all the subfiles are of equal size. This is a key property of the proposed placement scheme which allows us to efficiently transmit messages in the delivery phase which are simultaneously helpful for users who have requested files of different popularities.

Let us now turn to the delivery phase. To make matters concrete, let us suppose that the first three users have demanded  $W_1$  and the last user is interested in  $W_2$ . Therefore, our demand vector is  $d = (1, 1, 1, 2)$ . We define  $\Omega_i$  as the subset of users who have requested file  $W_i$ . In this case,  $\Omega_1 = \{1, 2, 3\}$  and  $\Omega_2 = \{4\}$ .

For the delivery phase, we construct a compressed description  $W_i^*$  for each file  $W_i$ . For users in  $\Omega_i$  recovering  $W_i^*$  implies recovering  $W_i$ , that is,  $H_q(W_i|W_i^*, Z_j) = 0$ . Moreover, among all  $W_i^*$  that satisfy this property, our particular construction minimizes both  $\max_{j \in \Omega_i} H_q(W_i^*|Z_j)$  and  $\max_{j \notin \Omega_i} H_q(W_i^*|Z_j)$  at the same time. The general construction of  $W_i^*$  is presented in Section VI, along with proofs

of its properties. For the example at hand, our construction specializes to

$$W_1^* = \left[ W_{1,\{4\},\{4\}}^*; W_{1,\{4\},\{1\}}^*; W_{1,\{1\},\{1\}}^* \right], \quad (8)$$

where

$$\begin{aligned} W_{1,\{4\},\{4\}}^* &= \left[ W_{1,\{4\},\{4\}}^{*(1)}; W_{1,\{4\},\{4\}}^{*(2)} \right], \\ W_{1,\{4\},\{1\}}^* &= \left[ W_{1,\{4\},\{1\}}^{*(1)}; W_{1,\{4\},\{1\}}^{*(2)} \right], \\ W_{1,\{1\},\{1\}}^* &= \left[ W_{1,\{1\},\{1\}}^{*(1)}; W_{1,\{1\},\{1\}}^{*(2)} \right], \end{aligned} \quad (9)$$

and

$$\begin{aligned} W_{1,\{4\},\{4\}}^{*(1)} &= W_{1,\{1,4\},\{4\}} + W_{1,\{2,4\},\{4\}} + W_{1,\{3,4\},\{4\}}, \\ W_{1,\{4\},\{4\}}^{*(2)} &= W_{1,\{1,4\},\{4\}} + 2W_{1,\{2,4\},\{4\}} + 3W_{1,\{3,4\},\{4\}}, \\ W_{1,\{4\},\{1\}}^{*(1)} &= W_{1,\{1,4\},\{1\}} + W_{1,\{2,4\},\{2\}} + W_{1,\{3,4\},\{3\}}, \\ W_{1,\{4\},\{1\}}^{*(2)} &= W_{1,\{1,4\},\{1\}} + 2W_{1,\{2,4\},\{2\}} + 3W_{1,\{3,4\},\{3\}}, \\ W_{1,\{1\},\{1\}}^{*(1)} &= W_{1,\{1,2\},\{1\}} + W_{1,\{1,2\},\{2\}} + W_{1,\{1,3\},\{1\}} \\ &\quad + W_{1,\{1,3\},\{3\}} + W_{1,\{2,3\},\{2\}} + W_{1,\{2,3\},\{3\}}, \\ W_{1,\{1\},\{1\}}^{*(2)} &= W_{1,\{1,2\},\{1\}} + 2W_{1,\{1,2\},\{2\}} + W_{1,\{1,3\},\{1\}} \\ &\quad + 2W_{1,\{1,3\},\{3\}} + W_{1,\{2,3\},\{2\}} + 2W_{1,\{2,3\},\{3\}}. \end{aligned} \quad (10)$$

Therefore, the subfiles of  $W_1^*$  are  $W_{1,\{4\},\{4\}}^{*(1)}, W_{1,\{4\},\{4\}}^{*(2)}, W_{1,\{4\},\{1\}}^{*(1)}, W_{1,\{4\},\{1\}}^{*(2)}, W_{1,\{1\},\{1\}}^{*(1)}, W_{1,\{1\},\{1\}}^{*(2)}$ . We can represent this in matrix form as follows.

$$W_1^* = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 2 & 0 & 3 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 2 & 0 & 3 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 2 & 1 & 2 & 0 & 0 & 1 & 2 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} W_{1,\{1,2\},\{1\}} \\ W_{1,\{1,2\},\{2\}} \\ W_{1,\{1,3\},\{1\}} \\ W_{1,\{1,3\},\{3\}} \\ W_{1,\{1,4\},\{1\}} \\ W_{1,\{1,4\},\{4\}} \\ W_{1,\{2,3\},\{2\}} \\ W_{1,\{2,3\},\{3\}} \\ W_{1,\{2,4\},\{2\}} \\ W_{1,\{2,4\},\{4\}} \\ W_{1,\{3,4\},\{3\}} \\ W_{1,\{3,4\},\{4\}} \end{bmatrix}. \quad (11)$$

If a user in  $\Omega_1$  successfully receives  $W_1^*$ , he can, with the help of side information already stored in his cache, recover the entire  $W_1$ . For instance, user 1 only needs to solve the following set of equations for  $W_{1,\{2,3\},\{2\}}, W_{1,\{2,3\},\{3\}}, W_{1,\{2,4\},\{2\}}, W_{1,\{2,4\},\{4\}}, W_{1,\{3,4\},\{3\}}, W_{1,\{3,4\},\{4\}}$ .

$$\begin{bmatrix} W_{1,\{4\},\{4\}}^{*(1)} \\ W_{1,\{4\},\{4\}}^{*(2)} \\ W_{1,\{4\},\{1\}}^{*(1)} \\ W_{1,\{4\},\{1\}}^{*(2)} \\ W_{1,\{1\},\{1\}}^{*(1)} \\ W_{1,\{1\},\{1\}}^{*(2)} \end{bmatrix} - \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 2 & 1 & 2 & 0 & 0 \end{bmatrix} \begin{bmatrix} W_{1,\{1,2\},\{1\}} \\ W_{1,\{1,2\},\{2\}} \\ W_{1,\{1,3\},\{1\}} \\ W_{1,\{1,3\},\{3\}} \\ W_{1,\{1,4\},\{1\}} \\ W_{1,\{1,4\},\{4\}} \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 2 & 0 & 3 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 2 & 0 & 3 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 2 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} W_{1,\{2,3\},\{2\}} \\ W_{1,\{2,3\},\{3\}} \\ W_{1,\{2,4\},\{2\}} \\ W_{1,\{2,4\},\{4\}} \\ W_{1,\{3,4\},\{3\}} \\ W_{1,\{3,4\},\{4\}} \end{bmatrix}. \quad (12)$$

This is possible since user 1 knows the left-hand side of the equation, and the matrix on the right hand-side is invertible. Therefore, our goal boils down to transferring the entire  $W_1^*$  to all the users in  $\Omega_1$ . Following a similar process, we construct the description  $W_2^*$  as follows

$$W_2^* = \left[ W_{2,\{1\},\{1\}}^*; W_{2,\{2\},\{2\}}^*; W_{2,\{3\},\{3\}}^*; W_{2,\{1,2\},\{1\}}^*; W_{2,\{1,2\},\{2\}}^*; W_{2,\{1,3\},\{1\}}^*; W_{2,\{1,3\},\{3\}}^*; W_{2,\{2,3\},\{2\}}^*; W_{2,\{2,3\},\{3\}}^* \right], \quad (13)$$

where

$$\begin{aligned} W_{2,\{i\},\{i\}}^* &= W_{2,\{i,4\},\{i\}} \quad \forall i \in \{1, 2, 3\}, \\ W_{2,\{i,j\},\{i\}}^* &= W_{2,\{i,j\},\{i\}} \quad \forall (i, j) \text{ s.t. } i, j \in [3], i \neq j. \end{aligned} \quad (14)$$

That is, in this example,  $W_2^*$  consists of the subfiles of  $W_2$  which are unknown to user 4. Again, transferring the entire  $W_2^*$  to user 4, guarantees his successful recovery of  $W_2$ .

To simplify matters, we will require every user in  $[K]$  to recover the entire  $[W_1^*; W_2^*]$ . In order to accomplish this, we transmit  $\mathcal{C}[W_1^*; W_2^*]$  over the broadcast link. The matrix  $\mathcal{C}$  here is an MDS matrix of 12 rows and 15 columns. The number of rows of this matrix is determined by the maximum number of subfiles of  $[W_1^*; W_2^*]$  which are unknown to any given user. In this example, a user in  $\Omega_1$  has precisely 12 unknowns in  $[W_1^*; W_2^*]$  (6 subfiles of  $W_1^*$  and 6 subfiles of  $W_2^*$ ). On the other hand, a user in  $\Omega_2$  knows 4 out of the 6 subfiles of  $W_1^*$ . Therefore, a total of 11 subfiles of  $[W_1^*; W_2^*]$  are unknown to him. Hence, the matrix  $\mathcal{C}$  must have 12 rows. Once a user in  $[K]$  receives  $\mathcal{C}[W_1^*; W_2^*]$ , he can remove the columns of  $\mathcal{C}$  which correspond to the subfiles he already knows. The resulting matrix will be square (or overdetermined) which is invertible owing to the MDS structure of  $\mathcal{C}$ . This will allow every user to decode  $[W_1^*; W_2^*]$ . Subsequently, each user in  $\Omega_i$  can proceed to decode  $W_i$  with the help of his side information. Recall that we started by dividing each file into 12 subfiles, and the delivery message consists of 12 linear combinations of such subfiles. Therefore, the delivery rate for this particular request vector is  $R = 1$ .

As we will see in the next section, the delivery rate of our strategy only depends on the request vector  $\mathbf{d}$  through  $\mathcal{N} = \text{Range}(\mathbf{d})$ , the set of indices of all the files that have been requested at least once. Therefore, we showed that with  $N = 2, K = 4, (r_1, r_2) = (2, 1)$ , and assuming  $\mathcal{N} = \{1, 2\}$ , we can achieve a delivery rate of  $R = 1$ . We can perform the same process for every choice of  $(r_1, r_2) \in \mathbb{Z}^2$  that satisfies  $0 \leq r_2 \leq r_1 \leq K$ . The result is summarized in Table II. Note that if  $\mathcal{N} = \{i\}$ , the delivery rate is simply  $1 - \frac{r_i}{K}$ .

$\begin{matrix} (r_1, r_2) \\ \mathcal{N} \end{matrix}$	(0, 0)	(1, 0)	(1, 1)	(2, 0)	(2, 1)
$\begin{matrix} \{1, 2\} \\ \{1\} \\ \{2\} \end{matrix}$	2 1 1	7/4 3/4 1	5/4 3/4 3/4	3/2 1/2 1	1 1/2 3/4
$\begin{matrix} (r_1, r_2) \\ \mathcal{N} \end{matrix}$	(2, 2)	(3, 0)	(3, 1)	(3, 2)	(3, 3)
$\begin{matrix} \{1, 2\} \\ \{1\} \\ \{2\} \end{matrix}$	2/3 1/2 1/2	5/4 1/4 1	3/4 1/4 3/4	1/2 1/4 1/2	1/4 1/4 1/4
$\begin{matrix} (r_1, r_2) \\ \mathcal{N} \end{matrix}$	(4, 0)	(4, 1)	(4, 2)	(4, 3)	(4, 4)
$\begin{matrix} \{1, 2\} \\ \{1\} \\ \{2\} \end{matrix}$	1 0 1	3/4 0 3/4	1/2 0 1/2	1/4 0 1/4	0 0 0

TABLE II: The set of delivery rates of our proposed scheme for all possible choices of  $0 \leq r_2 \leq r_1 \leq K$  and all possible  $\mathcal{N} \subseteq [N]$ . We have  $N = 2$  and  $K = 4$ .

By performing memory-sharing among all such points, we are able to achieve the lower convex envelope of the points in Table II. The expected delivery rate as a function of  $(r_1, r_2)$  for a probability distribution of  $(p_1, p_2) = (0.8, 0.2)$  has been plotted in Figure 2. Note that the dotted half of the figure where  $r_2 > r_1$  would correspond to switching the roles of the two files  $W_1$  and  $W_2$  and allocating more cache to the less popular file. The next question is how to find the best delivery rate for a particular cache size  $M$ . For this, we first have to restrict Figure 2 to the trajectory  $r_1 + r_2 = MK$ . As an example, we have plotted the thick red curve on the figure which corresponds to  $r_1 + r_2 = 3$  (or  $M = 3/4$ ). In order to find the best caching strategy for a cache size of  $M = 3/4$ , we need to choose the global minimum of this red curve. This can be done efficiently due to the convexity of the curve, and as we will see in Section VIII-A, can be even performed via binary search over the set of break points of the curve. As marked on the figure with a red circle, for this particular example with  $K = 4, N = 2, (p_1, p_2) = (0.8, 0.2), M = 3/4$ , the expected delivery rate is 0.79 which can be achieved by allocating a cache of size  $M_1 = r_1/K = 0.5$  to file  $W_1$  and  $M_2 = r_2/K = 0.25$  to file  $W_2$ . Theorem 3 from Section VIII will tell us that under the restriction of uncoded placement, this is the best expected delivery rate that one can achieve for the given  $(K, N, \mathbf{p}, M)$ .

## V. THE PLACEMENT STRATEGY

In this section we describe our general placement strategy. Note that our placement strategy can be applied to an arbitrary number of files and users and can even be adapted to user-specific file popularities (see Remark 2). Without loss of generality, suppose that the files are indexed in decreasing order of their popularity. In other words, file  $W_i$  is at least as popular as file  $W_{i+1}$  for all  $i \in [N - 1]$ . The placement strategy begins with selecting integers  $r_1, \dots, r_N$  such that  $0 \leq r_N \leq \dots \leq r_1 \leq K$ . Each  $r_i$  is proportional to the amount of cache that we are willing to allocate to file  $W_i$ . We divide each file into

$$S = \left( r_N, r_{N-1} - r_N, \dots, r_1 - r_2, K - r_1 \right) \quad (15)$$

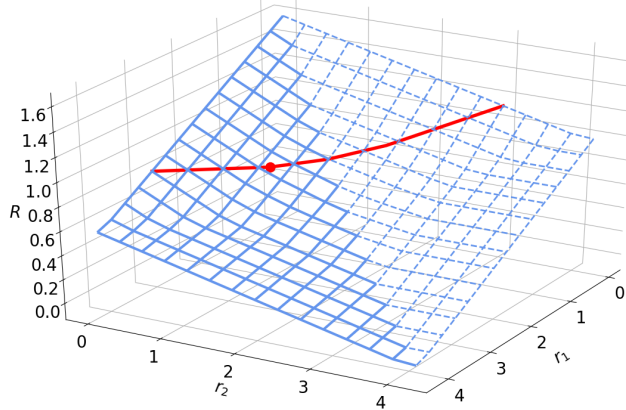


Fig. 2: The expected delivery rate for the caching problem with 4 users and 2 files versus  $(r_1, r_2)$ . The probabilities of the two files are 0.8 and 0.2 respectively. The thick red curve determines the set of  $(r_1, r_2)$  which results in a cache of size  $M = 3/4$ . The red circle on the curve is the minimizer of the red curve, which provides the optimal delivery rate under uncoded placement for the given  $(K, N, \mathbf{p}, M)$ .

subfiles of equal size. We label each subfile by  $N$  sets  $\tau_1, \dots, \tau_N$  where  $|\tau_j| = r_j$  for  $j \in [N]$  and  $\tau_j \subseteq \tau_{j-1}$  for  $j \in [2 : N]$  and  $\tau_1 \subseteq [K]$ . It should be evident that there are exactly  $S$  such subfiles. Next, for file  $W_i$ , we require each user  $k$  to store the subfile  $W_{i, \tau_1, \dots, \tau_N}$  if and only if  $k \in \tau_i$ . This process has been summarized in Algorithm 1, and an illustration for the case of  $N = 2$  has been provided in Figure 3. We can compute the amount of cache dedicated by each user to file  $i$  as follows

$$\begin{aligned} M_i &= \frac{\binom{K-r_2}{r_1-r_2} \times \dots \times \binom{K-r_i}{r_{i-1}-r_i} \binom{K-1}{r_i-1} \binom{r_i}{r_{i+1}} \times \dots \times \binom{r_{N-1}}{r_N}}{S} \\ &= \frac{\binom{K-r_2}{r_1-r_2} \times \dots \times \binom{K-r_i}{r_{i-1}-r_i} \binom{K-1}{r_i-1} \binom{r_i}{r_{i+1}} \times \dots \times \binom{r_{N-1}}{r_N}}{\binom{K-r_2}{r_1-r_2} \times \dots \times \binom{K-r_i}{r_{i-1}-r_i} \binom{K}{r_i} \binom{r_i}{r_{i+1}} \times \dots \times \binom{r_{N-1}}{r_N}} \\ &= \frac{r_i}{K}. \end{aligned} \quad (16)$$

This results in a total normalized cache size of

$$M = \sum_{i=1}^N M_i = \frac{\sum_{i=1}^N r_i}{K}. \quad (17)$$

*Remark 1.* In the special case of  $r_1 = \dots = r_N$ , all the sets  $\tau_i$  will be equal, and can be represented by only one set  $\tau$ . In this case, our placement phase is equivalent to the uniform placement strategy proposed in [1].

*Remark 2.* More generally, each user  $j$  could choose a permutation  $\pi : [N] \rightarrow [N]$  and store file  $W_{i, \tau_1, \dots, \tau_N}$  if and only if  $j \in \tau_{\pi(i)}$ . This would allow different users to have different preferences in terms of the popularities of the files, while still keeping all the cache sizes equal, and maintaining the same sub-packetization for all files. To provide a simple example, suppose we have two users and two files  $W_1$  and  $W_2$  with  $(r_1, r_2) = (2, 1)$  resulting in a cache size

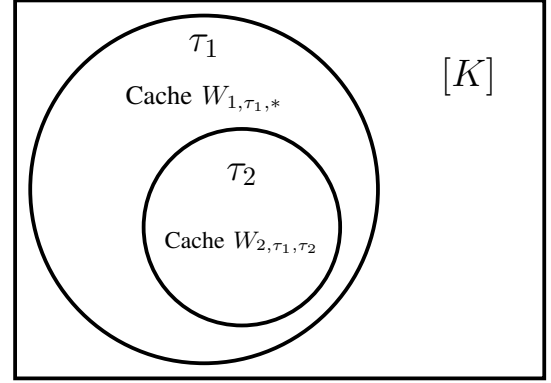


Fig. 3: Van Diagram of the placement strategy for the case of two files. Users whose indices appear in  $\tau_1$  cache  $W_{1, \tau_1, \tau_2}$  for all  $\tau_2 \subseteq \tau_1$ ,  $|\tau_2| = r_2$ . Users whose indices appear in  $\tau_2$  cache  $W_{2, \tau_1, \tau_2}$ .

of  $M = 3/2$ . In this case, user 1 could cache  $W_{1, \{1,2\}, \{1\}}$  and  $W_{1, \{1,2\}, \{2\}}$  but only  $W_{2, \{1,2\}, \{1\}}$ . On the other hand, user 2 could cache  $W_{2, \{1,2\}, \{1\}}$  and  $W_{2, \{1,2\}, \{2\}}$  but only  $W_{1, \{1,2\}, \{2\}}$ . This caching scheme preserves the property that each file has the same number of subfiles, while allowing each user to give higher priority to a different file.

---

#### Algorithm 1 The Placement Strategy for $N$ files and $K$ users

---

**Input:**  $(W_1, \dots, W_N), (r_1, \dots, r_N), K$

**Output:** The placement contents  $(Z_1, \dots, Z_K)$ .

- 1:  $S = \binom{K}{r_N, r_{N-1}-r_N, \dots, r_1-r_2, K-r_1}$ .
  - 2: Break each file  $W_i$  into  $S$  non-overlapping subfiles of equal size and index them as
 
$$W_i = \{W_{i, \tau_1, \dots, \tau_N} \mid \tau_N \subseteq \dots \subseteq \tau_1 \subseteq [K], |\tau_j| = r_j, \forall j \in [N]\}.$$
  - 3: **for**  $i \in [K]$  **do**
  - 4:    $Z_i = \{W_{j, \tau_1, \dots, \tau_N} \mid \text{for all } j \in [N] \text{ and all } (\tau_1, \dots, \tau_N) \text{ such that } i \in \tau_j\}.$
  - 5: **end for**
  - 6: Return  $(Z_1, \dots, Z_K)$ .
- 

*Remark 3.* The combinatorial designs in the Maddah-Ali and Niesen placement strategy [1] have recently been used in other closely related fields such as Coded Data Shuffling [25] and Coded Map-Reduce [26]. Despite being very useful at capturing the symmetric commonalities of different objects, it is not trivial how one can systematically generalize this combinatorial design to heterogeneous networks. We believe that our proposed multiset indexing tool provides a flexible, yet systematic extension of this scheme to asymmetric settings. For instance, consider a uniform caching problem with heterogeneous cache sizes. Let us assume that we have  $S$  different cache sizes  $M_1, \dots, M_S$ . Define  $(r_1, \dots, r_S) = (\frac{KM_1}{N}, \dots, \frac{KM_S}{N})$  and for simplicity assume  $r_i \in \mathbb{Z}$ . Divide each file  $W_i$  into subfiles of equal size  $W_{i, \tau_1, \dots, \tau_S}$  such that  $|\tau_j| = r_j$ ,  $\tau_j \subseteq \tau_{j-1}$  and  $\tau_1 \subseteq [K]$ . We could then let the

users with cache size  $M_j$  store  $W_{i,\tau_1,\dots,\tau_S}$  for all  $i \in [N]$  if and only if their indices are in  $\tau_j$ . As a second example, suppose we want to design a Fractional Repetition code [27] for distributed storage where some servers store more data than the others, and play a central role in the data-recovery criterion or the repair process in a distributed storage network. Again, one can systematically design such a storage code with the proposed multiset-indexing framework.

## VI. DELIVERY STRATEGY FOR $N = 2$

Let  $\Omega_1, \Omega_2 \subseteq [K]$  represent the subsets of the users that have requested files  $W_1$  and  $W_2$ , respectively. Therefore  $\Omega_1 \cap \Omega_2 = \emptyset$  and  $\Omega_1 \cup \Omega_2 = [K]$ . Also define  $K_1 = |\Omega_1|$  and  $K_2 = |\Omega_2| = K - K_1$ . Note that if  $K_1 = 0$ , a delivery rate of  $R = 1 - \frac{r_2}{K}$  can be trivially achieved. Similarly, if  $K_2 = 0$ , we can achieve a delivery rate of  $R = 1 - \frac{r_1}{K}$ . Let us now assume that both files have been requested.

The general idea behind the delivery scheme is as follows. First, we encode each file  $W_i, i \in [2]$  as  $W_i^*$  in such a way that decoding  $W_i^*$  provides enough information for each user in  $\Omega_i$  to decode  $W_i$ . In other words, we want that  $H_q(W_i|Z_j, W_i^*) = 0$  for all  $j \in \Omega_i, i \in [2]$ . The server transmits sufficient information for all the users in  $[K]$  to recover both  $W_1^*$  and  $W_2^*$ . Subsequently, each user in  $\Omega_i$  proceeds to recover  $W_i$  based on  $W_i^*$  and the contents of his cache. Moreover, the goal is for  $W_i^*$  to have a significant overlap with the cache of the users outside  $\Omega_i$ , i.e.,  $\max_{j \in [K] \setminus \Omega_i} H_q(W_i^*|Z_j)$  is as small as possible. The following lemma lays the foundation for our search for the ideal  $W_1^*$  and  $W_2^*$ .

**Lemma 1.** *Suppose  $\Omega_1 \neq \emptyset$  and  $\Omega_2 \neq \emptyset$ . For each  $i \in [2]$ , assume  $W_i^*$  satisfies*

$$H_q(W_i|W_i^*, Z_j) = 0, \quad \forall j \in \Omega_i. \quad (18)$$

Then  $W_i^*$  must satisfy

$$\max_{m \in \Omega_i} H_q(W_i^*|Z_m) \geq \frac{\binom{K-1}{r_i}}{\binom{K}{r_i}} F \quad (19)$$

and

$$\max_{\ell \in [K] \setminus \Omega_i} H_q(W_i^*|Z_\ell) \geq \frac{\binom{K-2}{r_i}}{\binom{K}{r_i}} F. \quad (20)$$

Furthermore, there exist  $W_1^*$  and  $W_2^*$  that satisfy Equations (18), (19) and (20) with equality.

*Proof of converse:* To prove Equation (19), note that

$$H_q(W_i^*|Z_m) = H_q(W_i^*, Z_m|Z_m) \geq H_q(W_i|Z_m). \quad (21)$$

But,  $H_q(W_i|Z_m)$  is the number of subfiles of  $W_i$  unknown to user  $m$ , multiplied by the size of one subfile, which is given by

$$\begin{aligned} H_q(W_i|Z_m) &= \frac{\prod_{j=1}^{i-1} \binom{K-r_{j+1}}{r_j-r_{j+1}} \binom{K-1}{r_i} \prod_{j=i+1}^N \binom{r_{j-1}}{r_j} F}{\prod_{j=1}^{i-1} \binom{K-r_{j+1}}{r_j-r_{j+1}} \binom{K}{r_i} \prod_{j=i+1}^N \binom{r_{j-1}}{r_j}} \\ &= \frac{\binom{K-1}{r_i}}{\binom{K}{r_i}} F. \end{aligned} \quad (22)$$

To prove Equation (20), let us concentrate on one arbitrary pair  $(m, \ell)$  where  $m \in \Omega_i$  and  $\ell \in [K] \setminus \Omega_i$ .

$$\begin{aligned} H_q(W_i^*|Z_\ell) &\geq H_q(W_i^*|Z_\ell, Z_m) = H_q(W_i^*, Z_m|Z_\ell, Z_m) \\ &\stackrel{(a)}{\geq} H_q(W_i|Z_\ell, Z_m) \stackrel{(b)}{=} \frac{\binom{K-2}{r_i}}{\binom{K}{r_i}} F, \end{aligned} \quad (23)$$

where (a) follows from Equation (18) and (b) is due to our placement strategy. To see why (b) holds, note that  $H_q(W_i|Z_\ell, Z_m)$  is the number of subfiles of  $W_i$  unknown to both users  $\ell$  and  $m$ , multiplied by the size of one subfile. This is given by

$$\begin{aligned} H_q(W_i|Z_\ell, Z_m) &= \frac{\prod_{j=1}^{i-1} \binom{K-r_{j+1}}{r_j-r_{j+1}} \binom{K-2}{r_i} \prod_{j=i+1}^N \binom{r_{j-1}}{r_j} F}{\prod_{j=1}^{i-1} \binom{K-r_{j+1}}{r_j-r_{j+1}} \binom{K}{r_i} \prod_{j=i+1}^N \binom{r_{j-1}}{r_j}} \\ &= \frac{\binom{K-2}{r_i}}{\binom{K}{r_i}} F. \end{aligned} \quad (24)$$

Most of this section will be dedicated to constructing  $W_1^*$  and  $W_2^*$  that satisfy the achievability part of Lemma 1. Once we have designed such  $W_1^*$  and  $W_2^*$ , we will construct a delivery message that helps all the users decode both.

Let  $(\rho_1, \rho_2)$  be an arbitrary pair of sets such that  $\rho_2 \subseteq \rho_1 \subseteq \Omega_1$  and  $s_i \triangleq |\rho_i| \leq r_i$  for  $i \in [2]$ . Let  $W_{2,\rho_1,\rho_2}$  be a column vector whose elements are the subfiles of  $W_2$  of the form  $W_{2,\rho_1 \cup x_1, \rho_2 \cup x_2}$  for all  $x_2 \subseteq x_1 \subseteq \Omega_2$ . The order of the elements in  $W_{2,\rho_1,\rho_2}$  is immaterial, as long as it is known to the users. This vector has  $\kappa_2(s_1, s_2)$  elements where

$$\kappa_2(s_1, s_2) = \binom{K_2}{r_2 - s_2} \binom{K_2 - (r_2 - s_2)}{r_1 - s_1 - (r_2 - s_2)}. \quad (25)$$

However, note that each user in  $\Omega_2$  knows all but  $\theta_2(s_1, s_2)$  elements of  $W_{2,\rho_1,\rho_2}$  where

$$\begin{aligned} \theta_2(s_1, s_2) &= \binom{K_2 - 1}{r_2 - s_2} \binom{K_2 - (r_2 - s_2)}{r_1 - s_1 - (r_2 - s_2)} \\ &= \frac{K_2 - (r_2 - s_2)}{K_2} \kappa_2(s_1, s_2). \end{aligned} \quad (26)$$

From the perspective of a user in  $\Omega_1$ , the story is entirely different. He either knows the entire  $W_{2,\rho_1,\rho_2}$  (if his index is in the set  $\rho_2$ ) or he does not know anything about  $W_{2,\rho_1,\rho_2}$ . We shall encode the vector  $W_{2,\rho_1,\rho_2}$  of length  $\kappa_2(s_1, s_2)$  as a new vector  $W_{2,\rho_1,\rho_2}^*$  of length  $\theta_2(s_1, s_2)$  in such a way that decoding  $W_{2,\rho_1,\rho_2}^*$  enables each user in  $\Omega_2$  to decode  $W_{2,\rho_1,\rho_2}$ . By doing so, we are simultaneously *aligning* the subfiles of  $W_2$  which are unknown to the users in  $\Omega_1$  to the extent possible.

Let  $\mathcal{C}_{2,s_1,s_2}$  be an arbitrary MDS matrix with  $\theta_2(s_1, s_2)$  rows and  $\kappa_2(s_1, s_2)$  columns. We know that if we remove any of  $\kappa_2(s_1, s_2) - \theta_2(s_1, s_2)$  columns of  $\mathcal{C}_{2,s_1,s_2}$ , the resulting square matrix is invertible. Define

$$W_{2,\rho_1,\rho_2}^* = \mathcal{C}_{2,s_1,s_2} W_{2,\rho_1,\rho_2}. \quad (27)$$

Let  $W_2^*$  be a vertical concatenation of all the vectors  $W_{\rho_1, \rho_2}^*$  for all  $(\rho_1, \rho_2)$ . Let us calculate the length of the vector  $W_2^*$ .

$$\begin{aligned}
\text{length}(W_2^*) &= \sum_{\rho_2 \subseteq \rho_1 \subseteq \Omega_1} \theta_2(|\rho_1|, |\rho_2|) \\
&= \sum_{s_1, s_2} \binom{K_1}{s_2} \binom{K_1 - s_2}{s_1 - s_2} \theta_2(s_1, s_2) \\
&= \sum_{s_1, s_2} \binom{K_1}{s_2} \binom{K_2 - 1}{r_2 - s_2} \binom{K_1 - s_2}{s_1 - s_2} \binom{K_2 - (r_2 - s_2)}{(r_1 - r_2) - (s_1 - s_2)} \\
&= \sum_{s_2} \binom{K_1}{s_2} \binom{K_2 - 1}{r_2 - s_2} \sum_{s_3} \binom{K_1 - s_2}{s_3} \binom{K_2 - (r_2 - s_2)}{(r_1 - r_2) - s_3} \\
&\stackrel{(a)}{=} \binom{K - 1}{r_2} \binom{K - r_2}{r_1 - r_2} = S \frac{\binom{K-1}{r_2} \binom{K-r_2}{r_1-r_2}}{\binom{K}{r_2} \binom{K-r_2}{r_1-r_2}} = S \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}}. \tag{28}
\end{aligned}$$

where we have defined  $s_3 = s_1 - s_2$ , and (a) follows from applying the Vandermonde identity to each summation. We can also compute the number of subfiles in  $W_2^*$  which are unknown to a user  $j \in \Omega_1$  as

$$\begin{aligned}
e_j &= \sum_{\substack{\rho_2 \subseteq \rho_1 \subseteq \Omega_1 \\ j \notin \rho_2}} \theta_2(|\rho_1|, |\rho_2|) \\
&= \sum_{\rho_2 \subseteq \rho_1 \subseteq \Omega_1} \binom{K_1 - 1}{s_2} \binom{K_1 - s_2}{s_1 - s_2} \theta_2(s_1, s_2) \\
&= \sum_{s_1, s_2} \binom{K_1 - 1}{s_2} \binom{K_2 - 1}{r_2 - s_2} \binom{K_1 - s_2}{s_1 - s_2} \binom{K_2 - r_2 + s_2}{r_1 - r_2 - s_1 + s_2} \\
&= \sum_{s_2} \binom{K_1 - 1}{s_2} \binom{K_2 - 1}{r_2 - s_2} \sum_{s_3} \binom{K_1 - s_2}{s_3} \binom{K_2 - r_2 + s_2}{r_1 - r_2 - s_3} \\
&= \binom{K - 2}{r_2} \binom{K - r_2}{r_1 - r_2} = S \frac{\binom{K-2}{r_2} \binom{K-r_2}{r_1-r_2}}{\binom{K}{r_2} \binom{K-r_2}{r_1-r_2}} = S \frac{\binom{K-2}{r_2}}{\binom{K}{r_2}}. \tag{29}
\end{aligned}$$

It is not difficult to see that  $H_q(W_2^*|Z_j) = e_j \frac{F}{S}$  which matches the lower-bound presented in Lemma 1, Equation (20). Furthermore,  $H_q(W_2^*|Z_m)$  for  $m \in \Omega_2$  is upper-bounded by  $\text{length}(W_2^*) \frac{F}{S}$  which matches Equation (19). Quite similarly, by reversing the roles of the two files in the description above, one can find a column vector  $W_1^*$  of length  $S \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}}$  with  $S \frac{\binom{K-2}{r_1}}{\binom{K}{r_1}}$  subfiles unknown to any user in  $\Omega_2$ . This also serves as a proof for the achievability part of Lemma 1. Now that we defined  $W_1^*$  and  $W_2^*$ , it is left to construct a delivery message that enables all the users in  $[K]$  to decode both. To accomplish this, we construct the delivery message as

$$X_d = \mathcal{C}[W_1^*; W_2^*], \tag{30}$$

where  $\mathcal{C}$  represent an MDS matrix with  $S \max \left\{ \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}} + \frac{\binom{K-2}{r_1}}{\binom{K}{r_1}}, \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-2}{r_2}}{\binom{K}{r_2}} \right\}$  rows and  $S \left( \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}} \right)$  columns. Note that the number of rows of  $\mathcal{C}$  is chosen to be the maximum number of subfiles of

$[W_1^*; W_2^*]$  unknown to any user in  $[K]$ . The resulting delivery rate is

$$R = \max \left\{ \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}} + \frac{\binom{K-2}{r_1}}{\binom{K}{r_1}}, \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-2}{r_2}}{\binom{K}{r_2}} \right\}. \tag{31}$$

The delivery strategy has been summarized in Algorithm 2.

---

**Algorithm 2** The Delivery Strategy for  $N = 2$  files and  $K$  users

---

**Input:**  $(W_1, W_2), (r_1, r_2), \mathbf{d}, K$

**Output:** The delivery message  $X_d$ .

- 1:  $\Omega_1 = \{i \in [K] | d_i = W_1\}$  and  $\Omega_2 = [K] \setminus \Omega_1$ .
  - 2:  $K_i = |\Omega_i|$  for  $i \in \{1, 2\}$ .
  - 3: **for**  $i \in [2]$  **do**
  - 4:   **if**  $\Omega_i = [K]$  **then**
  - 5:      $R = \frac{K-r_i}{K}$ .
  - 6:     Let  $\mathcal{C}$  be an  $SR$  by  $S$  MDS matrix.
  - 7:     Return  $X_d = \mathcal{C}W_i$ .
  - 8:   **end if**
  - 9: **end for**
  - 10: **for**  $i \in [2]$  **do**
  - 11:   **for**  $s_1 \in [0 : \min\{r_1, |\Omega_{[2] \setminus \{i\}}|\}]$  **do**
  - 12:     **for**  $s_2 \in [0 : \min\{r_2, s_1\}]$  **do**
  - 13:        $\kappa_i(s_1, s_2) = \binom{K_i}{r_2 - s_2} \binom{K_i - (r_2 - s_2)}{r_1 - s_1 - (r_2 - s_2)}$ .
  - 14:        $\theta_i(s_1, s_2) = \frac{K_i - (r_i - s_i)}{K_i} \kappa_i(s_1, s_2)$ .
  - 15:       Let  $\mathcal{C}_{i, s_1, s_2}$  be a  $\theta_i(s_1, s_2)$  by  $\kappa_i(s_1, s_2)$  MDS matrix.
  - 16:       **for**  $\rho_1 \subseteq \Omega_{[2] \setminus \{i\}}$  s.t.  $|\rho_1| = s_1$  **do**
  - 17:         **for**  $\rho_2 \subseteq \rho_1$  s.t.  $|\rho_2| = s_2$  **do**
  - 18:         Let  $W_{i, \rho_1, \rho_2}$  be a vertical concatenation of the subfiles  $\{W_{i, \rho_1 \cup x_1, \rho_2 \cup x_2} | x_2 \subseteq x_1 \subseteq \Omega_i\}$ .
  - 19:          $W_{i, \rho_1, \rho_2}^* = \mathcal{C}_{i, s_1, s_2} W_{i, \rho_1, \rho_2}$ .
  - 20:         **end for**
  - 21:       **end for**
  - 22:     **end for**
  - 23:   **end for**
  - 24: **end for**
  - 25:  $R = \max \left\{ \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}} + \frac{\binom{K-2}{r_1}}{\binom{K}{r_1}}, \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-2}{r_2}}{\binom{K}{r_2}} \right\}$ .
  - 26: Let  $\mathcal{C}$  be a  $SR$  by  $S \left( \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}} \right)$  MDS matrix.
  - 27: **for**  $i \in [2]$  **do**
  - 28:   Let  $W_i^*$  be a vertical concatenation of the vectors  $\{W_{i, \rho_1, \rho_2}^* | \rho_2 \subseteq \rho_1 \subseteq \Omega_{[2] \setminus \{i\}}\}$ .
  - 29: **end for**
  - 30: Return  $X_d = \mathcal{C}[W_1^*; W_2^*]$ .
- 

### A. Correctness

In this section we prove the correctness of Algorithm 2 by establishing that upon receiving  $X_d$  every user will be able to recover his requested file. The decoding process for each user is done in two phases reminiscent of a peeling algorithm. In the first phase, each user decodes both  $W_1^*$  and  $W_2^*$ . In the second phase, each user  $i$  discards  $W_{\bar{d}_i}^*$  where  $\bar{d}_i$  is the



index of the file that has not been requested by user  $i$ . He then proceeds to decode  $W_{d_i}$  using only  $W_i^*$  and the side information stored in his cache.

*Decoding Step 1:* First let us show that after receiving  $X_{\mathbf{d}}$ , every user can recover the entire  $[W_1^*; W_2^*]$ . Remember that  $\mathcal{C}$  is an MDS matrix with  $SR$  rows and  $S \left( \binom{K-1}{r_1} + \binom{K-1}{r_2} \right)$  columns. Let  $\gamma_i$  represent the set of columns of  $\mathcal{C}$  corresponding to the elements of  $[W_1^*; W_2^*]$  which user  $i$  already knows from the side information available in his cache. Let  $\mathcal{C}_{\gamma_i}$  represent the submatrix of  $\mathcal{C}$  obtained by removing the columns indexed in  $\gamma_i$ . If this matrix is square (or overdetermined), user  $i$  will be able to invert it and recover  $[W_1^*; W_2^*]$ . Therefore, we need to prove that  $|\gamma_i| \geq S \left( \binom{K-1}{r_1} + \binom{K-1}{r_2} \right) - SR$ . This inequality directly follows from Equations (28), (29) and (31):

$$\begin{aligned} |\gamma_i| &= S \frac{\binom{K-1}{r_{\bar{d}_i}}}{\binom{K}{r_{\bar{d}_i}}} - S \frac{\binom{K-2}{r_{\bar{d}_i}}}{\binom{K}{r_{\bar{d}_i}}} \\ &\geq S \min \left\{ \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} - \frac{\binom{K-2}{r_1}}{\binom{K}{r_1}}, \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}} - \frac{\binom{K-2}{r_2}}{\binom{K}{r_2}} \right\} \\ &= S \left( \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}} \right) - SR. \end{aligned} \quad (32)$$

*Decoding Step 2:* Now we show that if a user in  $\Omega_1$  has the entire  $W_1^*$ , he can decode it for  $W_1$ . Similarly, if a user in  $\Omega_2$  has the entire  $W_2^*$ , he can recover  $W_2$ . It should however be noted that users in  $\Omega_{d_i}$  will not be able to recover  $W_{d_i}$ . The proof idea is very similar to Step 1. Remember that  $W_1^* = \{W_{1,\rho_1,\rho_2}^* | \rho_2 \subseteq \rho_1 \subseteq \Omega_2\}$ . We will show that once a user in  $\Omega_1$  has access to  $W_{1,\rho_1,\rho_2}^*$  he will be able to decode  $\{W_{1,\rho_1 \cup x_1, \rho_2 \cup x_2} | x_2 \subseteq x_1 \subseteq \Omega_1\}$  for all  $\rho_2 \subseteq \rho_1 \subseteq \Omega_2$ . To see why, note that  $W_{1,\rho_1,\rho_2}^* = \mathcal{C}_{1,s_1,s_2} W_{1,\rho_1,\rho_2}$  where  $s_i = |\rho_i|$ . The matrix  $\mathcal{C}_{1,s_1,s_2}$  is an MDS matrix with  $\theta_1(s_1, s_2)$  rows and  $\kappa_1(s_1, s_2)$  columns. The number of subfiles of  $W_{1,\rho_1,\rho_2}$  unknown to user  $i \in \Omega_1$  is equal to  $\theta_1(s_1, s_2)$ . User  $i$  can thus discard the remaining  $\kappa_1(s_1, s_2) - \theta_1(s_1, s_2)$  columns of  $\mathcal{C}_{1,s_1,s_2}$  and invert the resulting square matrix in order to recover his unknowns.

## B. Expected Achievable Rate

To summarize, we characterized the delivery strategy for every choice of  $(M_1, M_2)$  of the form  $M_i = \frac{r_i}{K}$  with  $r_i \in [0 : K]$ , and a non-trivial request vector. Two questions are left to be addressed. First, what if  $M_i K$  is not an integer, and second, for a fixed total cache size of  $M$ , what are the optimal values of  $M_1$  and  $M_2$ ? To answer the first question, we observe that the lower convex envelope of all the points  $(M_1, M_2, R)$  with  $M_i K \in \mathbb{Z}$  is achievable by simply performing memory-sharing among such points. If  $(M_1 K, M_2 K)$  is not a pair of integers, we rely on this memory-sharing strategy to find an achievability scheme. We postpone the second question to Section VIII-A, once we have a better understanding of the optimal memory-sharing strategy for arbitrary  $(M_1, M_2)$ . For now, we write the achievable expected delivery rate as the

minimum over all possible choices of  $(M_1, M_2)$  that satisfy  $M_1 + M_2 = M$ .

**Theorem 1.** *Consider the coded caching problem with 2 files  $W_1$  and  $W_2$ , and  $K$  users each equipped with a cache of size  $M$ . Denote the probability of requesting file  $W_i$  by  $p_i$  where  $p_1 + p_2 = 1$ . Then the following expected delivery rate is achievable.*

$$\begin{aligned} \bar{R}(M) &= \min_{\substack{t_1, t_2 \\ t_1 + t_2 = KM}} \left( \frac{K - t_1}{K} p_1^K + \frac{K - t_2}{K} p_2^K \right. \\ &\quad \left. + (1 - p_1^K - p_2^K) \mathcal{L}_{\mathbf{r} \rightarrow \mathbf{t}} \max((R_1(r_1, r_2), R_2(r_1, r_2))) \right) \end{aligned} \quad (33)$$

where

$$\begin{aligned} R_1(r_1, r_2) &= \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-2}{r_2}}{\binom{K}{r_2}}, \\ R_2(r_1, r_2) &= \frac{\binom{K-2}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}}. \end{aligned} \quad (34)$$

for all  $(r_1, r_2) \in \mathbb{Z}^2$  s.t.  $0 \leq r_1, r_2 \leq K$ .

*Proof:* Based on the proof of correctness of Algorithm 2, we know that if we allocate a cache of size  $M_i = r_i/K$  to file  $i$  for  $i \in [2]$ , and if both files are requested, then we can achieve a delivery rate of  $R_{r_1, r_2} = \max\{R_1(r_1, r_2), R_2(r_1, r_2)\}$ . If only file  $i$  has been requested, then we can easily achieve a delivery rate of  $\frac{K-r_i}{K}$ . By performing memory-sharing among all such points  $(r_1, r_2)$ , we are able to achieve the lower convex envelope of all the points  $((r_1, r_2), \bar{R}_{r_1, r_2})$  where  $\bar{R}_{r_1, r_2} = \frac{K-r_1}{K} p_1^K + \frac{K-r_2}{K} p_2^K + (1 - p_1^K - p_2^K) R_{r_1, r_2}$ , and  $0 \leq r_1, r_2 \leq K$  and  $r_i \in \mathbb{Z}$ . By restricting this lower convex envelope to the plane which yields a cache of size  $M$ , we can characterize the achievable expected rate,  $\bar{R}_{t_1, t_2}$ , for all  $(t_1, t_2)$  s.t.  $t_1/K + t_2/K = M$  and  $0 \leq t_1, t_2 \leq K$ ,  $(t_1, t_2) \in \mathbb{R}^2$ . We can then choose the pair  $(t_1, t_2)$  for which the rate  $\bar{R}_{t_1, t_2}$  is minimized. ■

*Remark 4.* In the special case of  $r_1 = r_2$ , our joint placement and delivery strategy achieves the same delivery rate as in [3]. This is because  $\frac{\binom{K-1}{r}}{\binom{K}{r}} + \frac{\binom{K-2}{r}}{\binom{K}{r}} = \frac{\binom{K}{r+1} - \binom{K-2}{r+1}}{\binom{K}{r}}$ . Although the two placement strategies become equivalent when  $r_1 = r_2$ , the delivery strategies remain distinct. In other words, Algorithm 2 offers an alternative delivery strategy for the uniform caching of two files, which can be of independent interest.

*Remark 5.* The key of our proposed delivery strategy is to construct, for each requested file  $W_i$ , a compressed description  $W_i^*$  that will be decoded by every user, even those who did not request  $W_i$ . One of the main results of our paper is that for the case of two users, this leads to optimal performance. Unfortunately, the same basic strategy fails to attain the lower bound developed below in Section VII. Nonetheless, we conjecture that this lower bound can indeed be attained using our general placement scheme together with an improved delivery strategy. However, the delivery strategy must fundamentally rely on the alignment of the undesired messages at each user, in a more subtle way than through the computation of  $W_i^*$ . In

particular, such alignments must occur across the subfiles of multiple undesired files, not just one.

## VII. CONVERSE BOUND FOR UNCODED PLACEMENT

In this section we will derive a converse bound for the expected delivery rate under uncoded placement for arbitrary  $K$ ,  $N$  and  $\mathbf{p}$ . For a request vector  $\mathbf{d}$ , we define  $\text{Range}(\mathbf{d})$  as the set of indices of the files that are requested at least once in  $\mathbf{d}$ . That is,  $\text{Range}(\mathbf{d}) = \{i \in [N] \mid \exists j \in [K] \text{ s.t. } d_j = i\}$ . To prove our converse bound, we will follow in the footsteps of Lemma 2 in [3] which provides an uncoded converse bound for uniform file popularities. To start with, fix a request vector  $\mathbf{d}$  and let  $\mathcal{N} = \text{Range}(\mathbf{d})$  and let  $u_i$  be the index of an arbitrary user such that  $d_{u_i} = W_i$ , and let  $\mathcal{U} = \{u_i \mid i \in \mathcal{N}\}$ .

The general idea behind the proof in [3] is to construct a virtual user whose cache contains a subset of the symbols stored by the users in  $\mathcal{U}$ . This is done in such a way that the virtual user can recover all the files  $\{W_i \mid i \in \mathcal{N}\}$  after receiving the delivery message  $X_{\mathbf{d}}$ . For instance, let us fix a bijective function  $\pi : [|\mathcal{N}|] \rightarrow \mathcal{N}$ . The virtual user can store the entire cache of user  $u_{\pi(1)}$ , but since  $(X_{\mathbf{d}}, Z_{u_{\pi(1)}})$  enables him to decode  $W_{\pi(1)}$ , he will only cache the symbols in  $Z_{u_{\pi(2)}}$  which do not belong to  $W_{\pi(1)}$ . Similarly, he can discard the symbols in  $Z_{u_{\pi(3)}}$  which belong to either  $W_{\pi(1)}$  or  $W_{\pi(2)}$ , and so on. The converse bound is simply  $H_q(\{W_i \mid i \in \mathcal{N}\} \mid Z)$  where  $Z$  is the cache of the virtual user.

This converse bound depends on the particular request vector  $\mathbf{d}$  and the choice of the "leaders"  $\mathcal{U}$ . To remove these dependencies, one can take the average of the converse bound over all possible request vectors that share the same  $\text{Range}(\mathbf{d}) = \mathcal{N}$  as well as all possible choices of the leaders. Finally, note that the converse bound also depends on  $\pi(\cdot)$  which indicates in which order different files indexed in  $\mathcal{N}$  are processed by the virtual user. We create one virtual user for each  $\pi$ . Since every such virtual user must be able to recover all the files  $\{W_i \mid i \in \mathcal{N}\}$ , the overall converse bound will be the maximum of the  $|\mathcal{N}|!$  bounds obtained in this fashion.

**Theorem 2.** *Consider the problem of coded caching with  $N$  files with probabilities  $(p_1, \dots, p_N)$  and  $K$  users such that each user has a cache of size  $M$ . The expected delivery rate under uncoded placement must satisfy*

$$\bar{R} \geq \min_{\substack{\mathbf{t} \\ \sum t_i = MK \\ 0 \leq t_i \leq K}} \left[ \sum_{\substack{\mathcal{N} \subseteq [N] \\ \text{Range}(\mathbf{d}) = \mathcal{N}}} \sum_{\mathbf{d} \in [N]^K} \prod_{i=1}^K p_{d_i} \max_{\pi: [|\mathcal{N}|] \rightarrow \mathcal{N}} R_{\pi}(\mathbf{t}, \mathcal{N}) \right] \quad (35)$$

where the maximum is taken over all bijections  $\pi : [|\mathcal{N}|] \rightarrow \mathcal{N}$ , and

$$R_{\pi}(\mathbf{t}, \mathcal{N}) = \sum_{i \in [|\mathcal{N}|]} \left[ (1 - t_{\pi(i)} + \lfloor t_{\pi(i)} \rfloor) \frac{\binom{K-i}{\lfloor t_{\pi(i)} \rfloor}}{\binom{K}{\lfloor t_{\pi(i)} \rfloor}} + (t_{\pi(i)} - \lfloor t_{\pi(i)} \rfloor) \frac{\binom{K-i}{\lfloor t_{\pi(i)} \rfloor + 1}}{\binom{K}{\lfloor t_{\pi(i)} \rfloor + 1}} \right]. \quad (36)$$

*Proof:* The proof closely follows that of Lemma 2 in [3] with a few minor but important differences. To start with, we assume that all the users together have dedicated a total (normalized) cache of size  $t_i$  to file  $W_i$  where  $0 \leq t_i \leq K$ . Since each user has a cache of size  $M$ , we must have  $\sum_{i=1}^N t_i = KM$ . As can be seen in the statement of Theorem 2, the converse bound for a particular request vector  $\mathbf{d}$ , only depends on  $\mathbf{d}$  through  $\mathcal{N}$ , the set of indices of the files that have been requested at least once. For a fixed request vector, we also define  $\Omega_i \subseteq [K], i \in \mathcal{N}$  as the set of indices of the users who have requested file  $W_i$ . For  $i \in \mathcal{N}$ , let  $u_i \in \Omega_i$  be the index of an arbitrary user who has requested file  $W_i$ , and let  $\mathcal{U} = \{u_i \mid i \in \mathcal{N}\}$ . Suppose an auxiliary user has access to the entire cache of user  $u_{\pi(i)}$  except for the symbols which belong to the files within  $\{W_{\pi(\ell)} \mid \ell \in [|\mathcal{N}|], \ell < i\}$ , for all  $i \in [|\mathcal{N}|]$ . Provided that this auxiliary user has received  $X_{\mathbf{d}}$ , he must be able to recover all the files within  $\{W_i \mid i \in \mathcal{N}\}$ . For this to be feasible, the delivery rate must satisfy [3]

$$R(\mathbf{t}, \mathcal{N}) \geq \frac{1}{F} \sum_{i \in [|\mathcal{N}|]} \sum_{j=1}^F \mathbf{1}(\mathcal{K}_{\pi(i),j} \cap \{u_{\pi(\ell)} \mid \ell \in [|\mathcal{N}|], \ell \leq i\} = \emptyset), \quad (37)$$

where  $\mathcal{K}_{\pi(i),j}$  represents the subset of the users that have cached the  $j$ 'th symbol of file  $W_{\pi(i)}$ . We take the average of the expression above over all request vectors  $\mathbf{d}$  that have the same  $\text{Range}(\mathbf{d}) = \mathcal{N}$  and over all possible choices of the set  $\mathcal{U}$ . We obtain

$$R(\mathbf{t}, \mathcal{N}) \geq \frac{1}{F} \sum_{i \in [|\mathcal{N}|]} \sum_{j=1}^F \frac{\binom{K-|\mathcal{K}_{\pi(i),j}|}{i}}{\binom{K}{i}}. \quad (38)$$

Similarly, we can build a new virtual user for every possible bijection  $\pi : [|\mathcal{N}|] \rightarrow \mathcal{N}$ . Each virtual user, gives us a new converse bound. Therefore, we have

$$R(\mathbf{t}, \mathcal{N}) \geq \frac{1}{F} \max_{\pi: [|\mathcal{N}|] \rightarrow \mathcal{N}} \sum_{i \in [|\mathcal{N}|]} \sum_{j=1}^F \frac{\binom{K-|\mathcal{K}_{\pi(i),j}|}{i}}{\binom{K}{i}}. \quad (39)$$

Let  $a_{n,i}$  represent the number of symbols of file  $W_i$  cached by exactly  $n$  users, normalized by  $F$ . We can write

$$R(\mathbf{t}, \mathcal{N}) \geq \max_{\pi: [|\mathcal{N}|] \rightarrow \mathcal{N}} \sum_{i \in [|\mathcal{N}|]} \sum_{n=0}^K \frac{\binom{K-n}{i}}{\binom{K}{i}} a_{n,\pi(i)} = \max_{\pi: [|\mathcal{N}|] \rightarrow \mathcal{N}} \sum_{i \in [|\mathcal{N}|]} \sum_{n=0}^K \frac{\binom{K-i}{n}}{\binom{K}{n}} a_{n,\pi(i)}. \quad (40)$$

For any  $i$ , consider the sequence  $c_{n,i} = \frac{\binom{K-i}{n}}{\binom{K}{n}}, n \in [0 : K]$  where  $\binom{a}{b} = 0$  if  $b > a$ . Let  $g_i : \mathbb{R} \rightarrow \mathbb{R}$  be the continuous

piecewise linear function whose corner points are  $c_{n,i}$ . In other words,

$$g_i(x) = \begin{cases} (1-x+\lfloor x \rfloor)c_{\lfloor x \rfloor,i} + (x-\lfloor x \rfloor)c_{\lfloor x \rfloor+1,i} & \text{if } \lfloor x \rfloor \in [0:K-1], \\ (1-x)c_{0,i} + xc_{1,i} & \text{if } \lfloor x \rfloor < 0, \\ (K-x)c_{K-1,i} + (x-K+1)c_{K,i} & \text{if } \lfloor x \rfloor > K-1. \end{cases} \quad (41)$$

Note that  $g_i(x)$  is a convex function for any  $i \in [\mathcal{N}]$ . Furthermore, the sequence  $a_{n,i}$ ,  $n \in [0:K]$  satisfies  $\sum_{n=0}^K a_{n,i} = 1$  and  $a_{n,i} \geq 0$ . Therefore, by Jensen's inequality we have

$$\sum_{n=0}^K a_{n,i}c_{n,i} = \sum_{n=0}^K a_{n,i}g_i(n) \geq g_i\left(\sum_{n=0}^K na_{n,i}\right). \quad (42)$$

But note that  $\sum_{n=0}^K na_{n,i} = t_i$ . As a result,

$$\begin{aligned} \sum_{n=0}^K a_{n,i}c_{n,i} &\geq g_i(t_i) = (1-t_i+\lfloor t_i \rfloor)c_{\lfloor t_i \rfloor,i} \\ &\quad + (t_i-\lfloor t_i \rfloor)c_{\lfloor t_i \rfloor+1,i} \\ &= (1-t_i+\lfloor t_i \rfloor)\frac{\binom{K-i}{\lfloor t_i \rfloor}}{\binom{K}{\lfloor t_i \rfloor}} + (t_i-\lfloor t_i \rfloor)\frac{\binom{K-i}{\lfloor t_i \rfloor+1}}{\binom{K}{\lfloor t_i \rfloor+1}}. \end{aligned} \quad (43)$$

Based on this, we can continue to bound Equation (40) as

$$\begin{aligned} R(\mathbf{t}, \mathcal{N}) &\geq \\ &\max_{\pi: [\mathcal{N}] \rightarrow \mathcal{N}} \sum_{i \in [\mathcal{N}]} \left[ (1-t_{\pi(i)}+\lfloor t_{\pi(i)} \rfloor)\frac{\binom{K-i}{\lfloor t_{\pi(i)} \rfloor}}{\binom{K}{\lfloor t_{\pi(i)} \rfloor}} \right. \\ &\quad \left. + (t_{\pi(i)}-\lfloor t_{\pi(i)} \rfloor)\frac{\binom{K-i}{\lfloor t_{\pi(i)} \rfloor+1}}{\binom{K}{\lfloor t_{\pi(i)} \rfloor+1}} \right]. \end{aligned} \quad (44)$$

Taking the expected value of this expression over all  $\mathcal{N}$  and the minimum of the resulting expression over all possible  $(t_1, \dots, t_N)$  provides the desired lower bound. ■

*Remark 6.* The minimization problem in Theorem 2 can be solved with standard convex optimization tools thanks to the fact that the right hand side of Equation (35) as well as the minimization constraints are convex in  $\mathbf{t}$ . To establish this fact, one only needs to show that

$$J(x) = (1-x+\lfloor x \rfloor)\frac{\binom{K-i}{\lfloor x \rfloor}}{\binom{K}{\lfloor x \rfloor}} + (x-\lfloor x \rfloor)\frac{\binom{K-i}{\lfloor x \rfloor+1}}{\binom{K}{\lfloor x \rfloor+1}} \quad (45)$$

is convex in  $x$ . But this expression is piece-wise linear in  $x$ . So, it is sufficient to prove that the slopes of the consecutive segments of  $J(x)$  increase by  $x$ . Or, in other words,

$$\frac{\binom{K-i}{\lfloor x \rfloor+1}}{\binom{K}{\lfloor x \rfloor+1}} - \frac{\binom{K-i}{\lfloor x \rfloor}}{\binom{K}{\lfloor x \rfloor}} \leq \frac{\binom{K-i}{\lfloor x \rfloor+2}}{\binom{K}{\lfloor x \rfloor+2}} - \frac{\binom{K-i}{\lfloor x \rfloor+1}}{\binom{K}{\lfloor x \rfloor+1}}. \quad (46)$$

This fact can be proven via elementary manipulations and is omitted for conciseness.

## VIII. OPTIMALITY RESULT FOR $N = 2$

In this section we prove that for the special case of  $N = 2$ , the converse bound provided by Equation (35) is tight. Our proof of optimality also sheds light on the points which contribute to the lower convex envelope at each  $(M_1, M_2)$  in Equation (33). As it turns out, it is always sufficient to look at the vicinity of the point  $(M_1, M_2)$ , and perform memory-sharing among points of the form  $(r_1, r_2)$  where  $r_i \in \{\lfloor M_i K \rfloor, \lceil M_i K \rceil\}$ . We start with a useful observation and then present a corollary of Theorem 1.

**Proposition 1.** *Suppose  $K, r_1, r_2$  are three positive integers such that  $0 \leq r_2 < r_1 \leq K$ . We have*

$$\frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-2}{r_2}}{\binom{K}{r_2}} \geq \frac{\binom{K-2}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}} \quad (47)$$

if and only if  $r_1 + r_2 \leq K$ .

*Proof:* Define  $A_1 = \frac{\binom{K-2}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-1}{r_2}}{\binom{K}{r_2}}$  and  $A_2 = \frac{\binom{K-1}{r_1}}{\binom{K}{r_1}} + \frac{\binom{K-2}{r_2}}{\binom{K}{r_2}}$ . Each  $A_i$  can be computed as

$$A_i = \frac{(K-1)(2K-r_1-r_2) - r_i(K-r_i)}{K(K-1)}. \quad (48)$$

Therefore,

$$\begin{aligned} A_2 - A_1 &= \frac{r_1(K-r_1) - r_2(K-r_2)}{K(K-1)} \\ &= \frac{(r_1-r_2)(K-r_1-r_2)}{K(K-1)}. \end{aligned} \quad (49)$$

Given that  $r_1 > r_2$ , we have  $A_2 - A_1 \geq 0$  if and only if  $r_1 + r_2 \leq K$ . ■

**Corollary 1.** *For the caching problem with  $K$  users, two files with probabilities  $p_1, p_2$  and cache size  $M$ , the following expected delivery rate is achievable for any  $(t_1, t_2) \in \mathbb{R}^2$  that satisfies  $t_1 + t_2 = MK$ .*

$$\begin{aligned} \bar{R}_{t_1, t_2} &= p_1^K \frac{K-t_1}{K} + p_2^K \frac{K-t_2}{K} \\ &\quad + (1-p_1^K - p_2^K) \max(R_1, R_2), \end{aligned} \quad (50)$$

where

$$\begin{aligned} R_1 &= (1+\lfloor t_1 \rfloor - t_1)\frac{\binom{K-1}{\lfloor t_1 \rfloor}}{\binom{K}{\lfloor t_1 \rfloor}} + (t_1 - \lfloor t_1 \rfloor)\frac{\binom{K-1}{\lfloor t_1 \rfloor+1}}{\binom{K}{\lfloor t_1 \rfloor+1}} \\ &\quad + (1+\lfloor t_2 \rfloor - t_2)\frac{\binom{K-2}{\lfloor t_2 \rfloor}}{\binom{K}{\lfloor t_2 \rfloor}} + (t_2 - \lfloor t_2 \rfloor)\frac{\binom{K-2}{\lfloor t_2 \rfloor+1}}{\binom{K}{\lfloor t_2 \rfloor+1}}, \end{aligned} \quad (51)$$

$$\begin{aligned} R_2 &= (1+\lfloor t_1 \rfloor - t_1)\frac{\binom{K-1}{\lfloor t_1 \rfloor}}{\binom{K}{\lfloor t_1 \rfloor}} + (t_1 - \lfloor t_1 \rfloor)\frac{\binom{K-1}{\lfloor t_1 \rfloor+1}}{\binom{K}{\lfloor t_1 \rfloor+1}} \\ &\quad + (1+\lfloor t_2 \rfloor - t_2)\frac{\binom{K-1}{\lfloor t_2 \rfloor}}{\binom{K}{\lfloor t_2 \rfloor}} + (t_2 - \lfloor t_2 \rfloor)\frac{\binom{K-1}{\lfloor t_2 \rfloor+1}}{\binom{K}{\lfloor t_2 \rfloor+1}}. \end{aligned} \quad (52)$$

*Proof:* We distinguish between two regimes.

**Regime 1.**  $t_1 - \lfloor t_1 \rfloor + t_2 - \lfloor t_2 \rfloor \geq 1$ . We will perform memory sharing between three points  $(r_1, r_2) \in T$  where

$T = \{([t_1], [t_2] + 1), ([t_1] + 1, [t_2]), ([t_1] + 1, [t_2] + 1)\}$ . The coefficients that we use for memory-sharing are respectively  $\theta_1, \theta_2, \theta_3$  where  $\theta_1 = 1 - t_1 + [t_1]$ ,  $\theta_2 = 1 - t_2 + [t_2]$  and  $\theta_3 = 1 - \theta_1 - \theta_2 = t_1 - [t_1] + t_2 - [t_2] - 1$ . The amount of cache dedicated to file  $W_1$  is thus  $\frac{[t_1]\theta_1 + (1 + [t_1])\theta_2 + ([t_1] + 1)\theta_3}{K} = \frac{t_1}{K}$ . Similarly, the amount of cache dedicated to file  $W_2$  is  $\frac{t_2}{K}$ . If only one file  $i$  is requested, we can trivially achieve  $R = \frac{K - t_i}{K}$ . Let us assume both files have been requested. There are two possibilities. Either  $[t_1] + [t_2] + 2 \leq K$  or  $[t_1] + [t_2] + 1 \geq K$ . In the first case, we know that Inequality (47) holds for every  $(r_1, r_2) \in T$ . As a result, the following delivery rate can be achieved

$$\begin{aligned}
R &= (1 - t_1 + [t_1]) \frac{\binom{K-1}{[t_1]}}{\binom{K}{[t_1]}} + (1 - t_1 + [t_1]) \frac{\binom{K-2}{[t_2] + 1}}{\binom{K}{[t_2] + 1}} \\
&+ (1 - t_2 + [t_2]) \frac{\binom{K-1}{[t_1] + 1}}{\binom{K}{[t_1] + 1}} + (1 - t_2 + [t_2]) \frac{\binom{K-2}{[t_2]}}{\binom{K}{[t_2]}} \\
&+ (t_1 - [t_1] + t_2 - [t_2] - 1) \frac{\binom{K-1}{[t_1] + 1}}{\binom{K}{[t_1] + 1}} \\
&+ (t_1 - [t_1] + t_2 - [t_2] - 1) \frac{\binom{K-2}{[t_2] + 1}}{\binom{K}{[t_2] + 1}} \\
&= (1 + [t_1] - t_1) \frac{\binom{K-1}{[t_1]}}{\binom{K}{[t_1]}} + (t_1 - [t_1]) \frac{\binom{K-1}{[t_1] + 1}}{\binom{K}{[t_1] + 1}} \\
&+ (1 + [t_2] - t_2) \frac{\binom{K-2}{[t_2]}}{\binom{K}{[t_2]}} + (t_2 - [t_2]) \frac{\binom{K-2}{[t_2] + 1}}{\binom{K}{[t_2] + 1}} \\
&= R_1.
\end{aligned} \tag{53}$$

In the second case, the direction of Inequality (47) is reversed for all  $(r_1, r_2) \in T$ . In this case, the following delivery rate can be achieved

$$\begin{aligned}
R &= (1 - t_1 + [t_1]) \frac{\binom{K-2}{[t_1]}}{\binom{K}{[t_1]}} + (1 - t_1 + [t_1]) \frac{\binom{K-1}{[t_2] + 1}}{\binom{K}{[t_2] + 1}} \\
&+ (1 - t_2 + [t_2]) \frac{\binom{K-2}{[t_1] + 1}}{\binom{K}{[t_1] + 1}} + (1 - t_2 + [t_2]) \frac{\binom{K-1}{[t_2]}}{\binom{K}{[t_2]}} \\
&+ (t_1 - [t_1] + t_2 - [t_2] - 1) \frac{\binom{K-2}{[t_1] + 1}}{\binom{K}{[t_1] + 1}} \\
&+ (t_1 - [t_1] + t_2 - [t_2] - 1) \frac{\binom{K-1}{[t_2] + 1}}{\binom{K}{[t_2] + 1}} \\
&= (1 + [t_1] - t_1) \frac{\binom{K-2}{[t_1]}}{\binom{K}{[t_1]}} + (t_1 - [t_1]) \frac{\binom{K-2}{[t_1] + 1}}{\binom{K}{[t_1] + 1}} \\
&+ (1 + [t_2] - t_2) \frac{\binom{K-1}{[t_2]}}{\binom{K}{[t_2]}} + (t_2 - [t_2]) \frac{\binom{K-1}{[t_2] + 1}}{\binom{K}{[t_2] + 1}} \\
&= R_2.
\end{aligned} \tag{54}$$

Therefore, we are able to achieve  $\max(R_1, R_2)$ . By taking the expectation over all possible request vectors, we achieve  $\bar{R}_{t_1, t_2}$  as in Equation (50).

**Regime 2.**  $t_1 - [t_1] + t_2 - [t_2] \leq 1$ . In this case we choose our set  $T = \{([t_1] + 1, [t_2]), ([t_1], [t_2] + 1), ([t_1], [t_2])\}$  and our coefficients  $\theta_1 = t_1 - [t_1]$ ,  $\theta_2 = t_2 - [t_2]$  and  $\theta_3 = 1 - t_1 + [t_1] - t_2 + [t_2]$ . Again we perform memory-sharing between the three points in  $T$  with the given coefficients. This allows us to achieve a delivery rate of  $\frac{K - t_i}{K}$  if only file  $W_i$  is requested. If both are requested, we can achieve  $R = \max\{R_1, R_2\}$ . By taking the expectation over all request vectors, we find the same delivery rate as in Equation (50). ■

It is easy to see that the achievable rate characterized by Corollary 1 lies on our converse bound in Equation (35). This implies that at any cache allocation point  $(t_1, t_2)$ , there are only three points  $(r_1, r_2)$  that contribute to the lower convex envelope. We first check whether  $t_1 - [t_1] + t_2 - [t_2] \geq 1$ . If this inequality holds, then we perform memory sharing among the three points  $\{([t_1], [t_2] + 1), ([t_1] + 1, [t_2]), ([t_1] + 1, [t_2] + 1)\}$ . Otherwise, we perform memory-sharing among  $\{([t_1] + 1, [t_2]), ([t_1], [t_2] + 1), ([t_1], [t_2])\}$ . Based on this observation, we can summarize our joint placement and delivery strategy as in Algorithm 3. The next theorem follows immediately from Corollary 1 and Theorem 2.

---

**Algorithm 3** The joint placement-delivery strategy for  $N = 2$  and arbitrary  $M, K, \mathbf{p}$

---

**Input:**  $W_1, W_2, M, K, \mathbf{p}$

**Output:** The cache contents  $(Z_1, \dots, Z_K)$  and all the delivery messages  $\{X_{\mathbf{d}} | \mathbf{d} \in [N]^K\}$ .

- 1:  $\mathbf{t} = (t_1, t_2) = \arg \min \bar{R}_{t_1, t_2}$  where  $\bar{R}_{t_1, t_2}$  is given by Equation (50).
- 2:  $r_1 = [t_1], r_2 = [t_2]$ .

**Placement**

- 3: **if**  $t_1 - r_1 + t_2 - r_2 \geq 1$  **then**
- 4:    $\theta_0 = 0, \theta_1 = 1 - t_1 + r_1, \theta_2 = 1 - t_2 + r_2,$   
     $\theta_3 = 1 - \theta_1 - \theta_2.$
- 5:    $Q_1 = (r_1, r_2 + 1), Q_2 = (r_1 + 1, r_2),$   
     $Q_3 = (r_1 + 1, r_2 + 1).$
- 6: **else**
- 7:    $\theta_0 = 0, \theta_1 = t_1 - r_1, \theta_2 = t_2 - r_2, \theta_3 = 1 - \theta_1 - \theta_2.$
- 8:    $Q_1 = (r_1 + 1, r_2), Q_2 = (r_1, r_2 + 1), Q_3 = (r_1, r_2).$
- 9: **end if**
- 10:  $P_j = \sum_{i=0}^j \theta_i$  for  $j \in [0 : 3]$ .
- 11:  $W_i^j = W_i |_{[P_{j-1}F+1:P_jF]}$  for  $i \in [2], j \in [3]$  where  $W_i |_A$  refers to the symbols of  $W_i$  indexed in the set  $A$ .
- 12:  $(Z_1^i, \dots, Z_K^i) =$  output of Algorithm 1 applied on  $((W_1^i, W_2^i), Q_i, K)$  for  $i \in [3]$ .
- 13:  $Z_j = (Z_j^1, Z_j^2, Z_j^3)$  for  $j \in [K]$ .

**Delivery**

- 14: **for** all request vectors  $\mathbf{d}$  **do**
  - 15:    $X_{\mathbf{d}}^{(i)} =$  output of Algorithm 2 applied on  $((W_1^i, W_2^i), Q_i, \mathbf{d}, K)$  for  $i \in [3]$ .
  - 16:    $X_{\mathbf{d}} = (X_{\mathbf{d}}^{(1)}, X_{\mathbf{d}}^{(2)}, X_{\mathbf{d}}^{(3)}).$
  - 17: **end for**
  - 18: **Return**  $((Z_1, \dots, Z_K), \{X_{\mathbf{d}} | \mathbf{d} \in [N]^K\})$ .
-

**Theorem 3.** For the coded caching problem with  $K$  users, two files with probabilities  $p_1, p_2$  and cache size  $M$ , the optimal expected delivery rate under uncoded placement is

$$\bar{R}^* = \min_{\substack{0 \leq t_2 \leq t_1 \leq K \\ t_1 + t_2 = KM}} \bar{R}_{t_1, t_2}, \quad (55)$$

where  $\bar{R}_{t_1, t_2}$  is given by Equation (50). Furthermore, this can be achieved by the joint placement and delivery strategy in Algorithm 3.

#### A. Finding the optimal memory allocation

The delivery rate in Equation (50) as a function of  $(t_1, t_2)$  is convex. As a result, the optimal  $(t_1, t_2)$ , which is the solution to

$$(t_1^*, t_2^*) = \arg \min_{\substack{0 \leq t_2 \leq t_1 \leq K \\ t_1 + t_2 = KM}} \bar{R}_{t_1, t_2}, \quad (56)$$

can be found by solving a convex optimization problem. However, note that our delivery rate is in fact a piece-wise linear function of  $(t_1, t_2)$ , the break points of which can be easily characterized. Based on the following theorem, we can find the optimal  $(t_1, t_2)$  by simply performing binary search over a discrete set of feasible points.

**Theorem 4.** There exists a solution  $(t_1^*, t_2^*)$  to Equation (56) that satisfies  $(t_1^*, t_2^*) \in \mathcal{P}$  and

$$\begin{aligned} m^+(t_1^*) &\geq \frac{p_1^K - (1-p_1)^K}{K(1-p_1^K - (1-p_1)^K)}, \\ m^-(t_1^*) &\leq \frac{p_1^K - (1-p_1)^K}{K(1-p_1^K - (1-p_1)^K)}, \end{aligned} \quad (57)$$

where we define

$$m^+(t_1) \triangleq \begin{cases} \left( \frac{\binom{K-1}{\lfloor t_1 \rfloor + 1}}{\binom{K-1}{\lfloor t_1 \rfloor}} - \frac{\binom{K-1}{\lfloor t_1 \rfloor}}{\binom{K-1}{\lfloor t_1 \rfloor}} + \frac{\binom{K-2}{\lceil KM-t_1 \rceil - 1}}{\binom{K-2}{\lceil KM-t_1 \rceil - 1}} - \frac{\binom{K-2}{\lceil KM-t_1 \rceil}}{\binom{K-2}{\lceil KM-t_1 \rceil}} \right) & \text{if } M \leq 1, \\ \left( \frac{\binom{K-2}{\lfloor t_1 \rfloor + 1}}{\binom{K-2}{\lfloor t_1 \rfloor}} - \frac{\binom{K-2}{\lfloor t_1 \rfloor}}{\binom{K-2}{\lfloor t_1 \rfloor}} + \frac{\binom{K-1}{\lceil KM-t_1 \rceil - 1}}{\binom{K-1}{\lceil KM-t_1 \rceil - 1}} - \frac{\binom{K-1}{\lceil KM-t_1 \rceil}}{\binom{K-1}{\lceil KM-t_1 \rceil}} \right) & \text{if } M > 1. \end{cases}$$

$$m^-(t_1) \triangleq \begin{cases} \left( \frac{\binom{K-1}{\lfloor t_1 \rfloor}}{\binom{K-1}{\lfloor t_1 \rfloor}} - \frac{\binom{K-1}{\lfloor t_1 \rfloor - 1}}{\binom{K-1}{\lfloor t_1 \rfloor - 1}} + \frac{\binom{K-2}{\lfloor KM-t_1 \rfloor}}{\binom{K-2}{\lfloor KM-t_1 \rfloor}} - \frac{\binom{K-2}{\lfloor KM-t_1 \rfloor + 1}}{\binom{K-2}{\lfloor KM-t_1 \rfloor + 1}} \right) & \text{if } t_1 > t_2 \text{ and } M \leq 1, \\ \left( \frac{\binom{K-2}{\lfloor t_1 \rfloor}}{\binom{K-2}{\lfloor t_1 \rfloor}} - \frac{\binom{K-2}{\lfloor t_1 \rfloor - 1}}{\binom{K-2}{\lfloor t_1 \rfloor - 1}} + \frac{\binom{K-1}{\lfloor KM-t_1 \rfloor}}{\binom{K-1}{\lfloor KM-t_1 \rfloor}} - \frac{\binom{K-1}{\lfloor KM-t_1 \rfloor + 1}}{\binom{K-1}{\lfloor KM-t_1 \rfloor + 1}} \right) & \text{if } t_1 > t_2 \text{ and } M > 1, \\ -m^+(t_1) & \text{if } t_1 = t_2. \end{cases}$$

$$\begin{aligned} \mathcal{P} &\triangleq \left\{ (t_1, t_2) \in \mathbb{R}^2 \mid 0 \leq t_2 \leq t_1 \leq K, \ t_1 + t_2 = KM, \right. \\ &\quad \left. (t_1 - \lfloor t_1 \rfloor)(t_2 - \lfloor t_2 \rfloor)(t_1 - t_2) = 0 \right\}. \end{aligned} \quad (58)$$

*Proof:* As we limit our achievable delivery rate to a line  $t_1 + t_2 = KM$ , we obtain a piecewise linear and convex curve  $\bar{R}(t_1)$ . It can be readily seen from Equation (50) that the break points of this curve occur when  $t_1 \in \mathbb{Z}$  or  $t_2 \in \mathbb{Z}$  or at the extreme point when  $t_1 = t_2$ . This establishes the choice of the feasible sets  $\mathcal{P}$  in the statement of the theorem.

If  $M \leq 1$ , we know that  $R_1$  is the maximizer of Equation (50). We can thus rephrase the expected delivery rate as

$$\begin{aligned} \bar{R}_{t_1, t_2} &= (1 - p_1^K - p_2^K) \left[ (1 + \lfloor t_1 \rfloor - t_1) \frac{\binom{K-1}{\lfloor t_1 \rfloor}}{\binom{K-1}{\lfloor t_1 \rfloor}} \right. \\ &\quad + (t_1 - \lfloor t_1 \rfloor) \frac{\binom{K-1}{\lfloor t_1 \rfloor + 1}}{\binom{K-1}{\lfloor t_1 \rfloor + 1}} + (1 + \lfloor t_2 \rfloor - t_2) \frac{\binom{K-2}{\lfloor t_2 \rfloor}}{\binom{K-2}{\lfloor t_2 \rfloor}} \\ &\quad \left. + (t_2 - \lfloor t_2 \rfloor) \frac{\binom{K-2}{\lfloor t_2 \rfloor + 1}}{\binom{K-2}{\lfloor t_2 \rfloor + 1}} \right] + p_1^K \frac{K - t_1}{K} + p_2^K \frac{K - t_2}{K} \\ &= t_1 \left[ (1 - p_1^K - (1 - p_1)^K) \left( \frac{\binom{K-1}{\lfloor t_1 \rfloor + 1}}{\binom{K-1}{\lfloor t_1 \rfloor + 1}} - \frac{\binom{K-1}{\lfloor t_1 \rfloor}}{\binom{K-1}{\lfloor t_1 \rfloor}} \right) \right. \\ &\quad \left. + \frac{\binom{K-2}{\lfloor KM-t_1 \rfloor}}{\binom{K-2}{\lfloor KM-t_1 \rfloor}} - \frac{\binom{K-2}{\lfloor KM-t_1 \rfloor + 1}}{\binom{K-2}{\lfloor KM-t_1 \rfloor + 1}} \right] + \frac{(1 - p_1)^K - p_1^K}{K} + c. \end{aligned} \quad (59)$$

Define  $\mathcal{Q} \triangleq \{t_1 \in \mathbb{R} \mid \exists t_2 \in \mathbb{R} \text{ s.t. } (t_1, t_2) \in \mathcal{P}\}$ . As long as  $t_1$  is in the open interval between two fixed consecutive members of  $\mathcal{Q}$ , the value of  $c$  does not change. As a result, the expression above provides us with the slope of the line segment which connects two consecutive points in the piecewise linear function  $\bar{R}(t_1)$ . Our goal is to find the value of  $t_1 \in \mathcal{Q}$  such that the slope of this curve is non-negative at  $t_1 + \epsilon$  and non-positive at  $t_1 - \epsilon$ . This is given by Equation (57). Note that we are using the identity  $(1 + \lfloor a \rfloor - a)g(\lfloor a \rfloor) + (a - \lfloor a \rfloor)g(\lfloor a \rfloor + 1) = (\lceil a \rceil - a)g(\lceil a \rceil - 1) + (a - \lceil a \rceil + 1)g(\lceil a \rceil)$  to simplify the expressions for  $m^+(t_1)$  and  $m^-(t_1)$ . Similar analysis can be made if  $M > 1$ . ■

## IX. NUMERICAL RESULTS

In this section we provide a numerical analysis of our caching strategy and compare it with the literature. First, we fix  $K = 6$ ,  $M = 1$  and  $N = 2$ , and find the optimal expected delivery rate (under uncoded placement) as a function of the probability of the first file, using Theorem 3. To accomplish this, we first have to find the optimal  $(t_1, t_2) = (M_1 K, M_2 K)$  as a function of  $p_1$  following Theorem 4. This optimal expected delivery rate has been plotted in Figure 4. A few valuable insights can be gained from this curve. Firstly, when the probabilities of the two files are close, the heuristic approach of applying uniform coded caching is indeed optimal. The range of probabilities for which this property holds ultimately depends on  $K$  and  $M$ , but for our example is given by  $|p_1 - p_2| < 0.48$ . This is the region that has been marked by  $(t_1, t_2) = (3, 3)$  on the figure. Similarly, when one file is very popular (in this case  $|p_1 - p_2| > 0.78$ ), it is optimal to allocate the entire cache to it, and ignore the other file in the placement phase. This region has been labeled as  $(t_1, t_2) = (6, 0)$ .

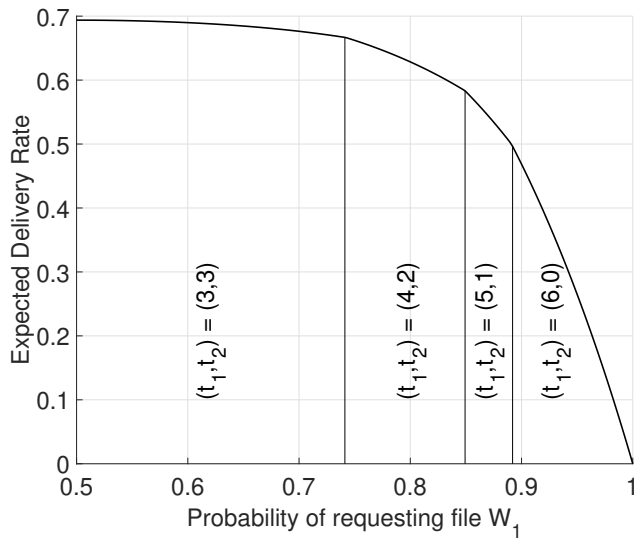


Fig. 4: The optimal expected delivery rate (under uncoded placement) for the non-uniform caching problem with  $K = 6$ ,  $N = 2$  and  $M = 1$ , versus the probability of requesting  $W_1$ .

But perhaps the most interesting scenario is when the probabilities lie somewhere in between. Figure 4 tells us that there is a range of probabilities (in this example  $0.48 < |p_1 - p_2| < 0.78$ ) for which no memory-sharing strategy is optimal. For this range, one must rely on Algorithm 3 with non-trivial choices of  $(t_1, t_2)$  to attain the optimal expected delivery rate. For this particular example, we must set  $(t_1, t_2) = (4, 2)$  for  $0.48 < |p_1 - p_2| < 0.70$  and  $(t_1, t_2) = (5, 1)$  for  $0.70 < |p_1 - p_2| < 0.78$ .

Now, let us instead fix  $p_1$  and find the optimal expected delivery rate as a function of the cache size  $M$ . For  $K = 6$ , we rely on the previous plot to choose  $p_1 = 0.85$  in order to emphasize the scenario where grouping is strictly sub-optimal, at least at  $M = 1$ . In Figure 5 we have plotted the optimal expected delivery rate for this choice of  $p_1$ , and compared it to the two possible grouping strategies [15]–[18]:  $R_{UN}$  corresponds to the uniform caching which ignores the differences in the probabilities of the two files, whereas  $R_{NC}$  is the delivery rate for a caching strategy that creates two groups each containing one file, and ignores the coding opportunities between the two. The expression for  $R_{UN}$  can be given [3] by the lower convex envelop of the points

$$R_{UN} = p_1^K \left(1 - \frac{r}{K}\right) + (1 - p_1)^K \left(1 - \frac{r}{K}\right) + (1 - p_1^K - (1 - p_1)^K) \left[ \frac{\binom{K}{r+1} - \binom{K-2}{r+1}}{\binom{K}{r}} \right], \quad (60)$$

where  $r = \frac{KM}{2} \in \mathbb{N}$ . As for  $R_{NC}$ , it is easy to see that if  $M \leq 1$ , the best memory allocation is to assign the entire cache to  $W_1$ . If  $M > 1$ , the remaining memory is given to file  $W_2$ . This results in a delivery rate of

$$R_{NC} = \begin{cases} (1 - p_1^K)(2 - M) & \text{if } M > 1, \\ -p_1^K - (1 - p_1)^K(1 - M) + 2 - M & \text{if } M \leq 1. \end{cases} \quad (61)$$

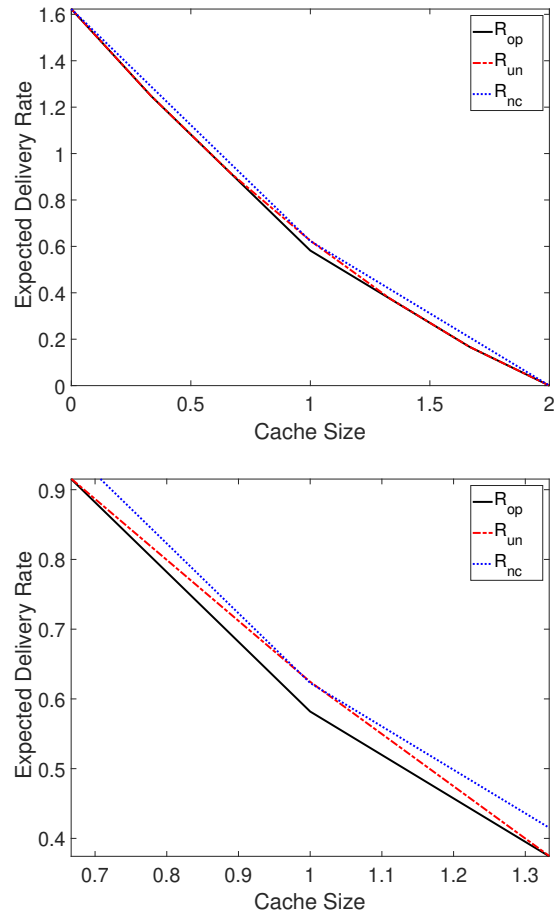


Fig. 5: **top**: Comparison of the expected delivery rate for our scheme  $R_{Op}$  (optimal under uncoded placement), and the two possible grouping strategies:  $R_{UN}$  performs uniform caching and ignores the differences in the probabilities, whereas  $R_{NC}$  ignores the coding opportunities between the two files. The parameters are  $K = 6$ ,  $N = 2$  and  $p_1 = 0.85$ . **bottom**: zoomed in on the vicinity of  $M = 1$ .

As visible in Figure 5, the most discrepancy between the grouping and optimal strategies occur around  $M = 1$ , where the optimal expected delivery rate is 0.582, whereas  $R_{nc} \approx 0.623$  and  $R_{un} \approx 0.625$ , about 7 percent larger than the optimal rate. On the other hand, at the extreme values of  $M$ , all three strategies are optimal. This is not very surprising: if  $M = 0$  or  $M = 2$ , all three strategies are equivalent. It is therefore natural that in the vicinity of such extreme values there is no major difference in their performances.

## X. CONCLUDING REMARKS

The majority of the existing literature on coded caching with non-uniform demands is focused on grouping strategies which can achieve constant additive or multiplicative gaps to the optimal expected delivery rate. This paper serves as a step towards the ambitious task of designing nonuniform coded caching strategies which are *optimal* under uncoded placement. Moreover, we believe that there is great potential to the multiset indexing extension of the uniform placement

strategy proposed in this paper, as it can be readily applied to other combinatorial problems of heterogeneous nature. Our delivery strategy for the case of two files may also serve as a stepping stone for a closer investigation of the application of interference alignment in coded caching.

#### ACKNOWLEDGEMENT

The authors would like to thank the Associate Editor and the reviewers for their invaluable feedback.

#### REFERENCES

- [1] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Transactions on Information Theory*, vol. 60, no. 5, pp. 2856–2867, 2014.
- [2] Q. Yu, M. A. Maddah-Ali, and A. S. Avestimehr, "Characterizing the rate-memory tradeoff in cache networks within a factor of 2," in *IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2017, pp. 386–390.
- [3] —, "The exact rate-memory tradeoff for caching with uncoded prefetching," *IEEE Transactions on Information Theory*, vol. 64, no. 2, pp. 1281–1296, 2018.
- [4] R. Pedarsani, M. A. Maddah-Ali, and U. Niesen, "Online coded caching," *IEEE/ACM Transactions on Networking (TON)*, vol. 24, no. 2, pp. 836–845, 2016.
- [5] N. Karamchandani, U. Niesen, M. A. Maddah-Ali, and S. N. Diggavi, "Hierarchical coded caching," *IEEE Transactions on Information Theory*, vol. 62, no. 6, pp. 3212–3229, 2016.
- [6] M. Ji, G. Caire, and A. F. Molisch, "Fundamental limits of caching in wireless d2d networks," *IEEE Transactions on Information Theory*, vol. 62, no. 2, pp. 849–869, 2016.
- [7] H. Ghasemi and A. Ramamoorthy, "Improved lower bounds for coded caching," *IEEE Transactions on Information Theory*, vol. 63, no. 7, pp. 4388–4413, 2017.
- [8] N. Golrezaei, K. Shanmugam, A. G. Dimakis, A. F. Molisch, and G. Caire, "Femtocaching: Wireless video content delivery through distributed caching helpers," in *INFOCOM, Proceedings IEEE*. IEEE, 2012, pp. 1107–1115.
- [9] C.-Y. Wang, S. H. Lim, and M. Gastpar, "A new converse bound for coded caching," in *Information Theory and Applications Workshop (ITA), 2016*. IEEE, 2016, pp. 1–6.
- [10] M. A. Maddah-Ali and U. Niesen, "Decentralized coded caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Transactions On Networking*, vol. 23, no. 4, pp. 1029–1040, 2015.
- [11] J. Gómez-Vilardebó, "Fundamental limits of caching: Improved bounds with coded prefetching," *arXiv preprint arXiv:1612.09071*, 2016.
- [12] K. Wan, D. Tuninetti, and P. Piantanida, "On the optimality of uncoded cache placement," in *IEEE Information Theory Workshop (ITW)*. IEEE, 2016, pp. 161–165.
- [13] Z. Chen, P. Fan, and K. B. Letaief, "Fundamental limits of caching: Improved bounds for small buffer users," *arXiv preprint arXiv:1407.1935*, 2014.
- [14] M. M. Amiri and D. Gunduz, "Fundamental limits of coded caching: Improved delivery rate-cache capacity trade-off," *arXiv preprint arXiv:1604.03888*, 2016.
- [15] U. Niesen and M. A. Maddah-Ali, "Coded caching with nonuniform demands," *IEEE Transactions on Information Theory*, vol. 63, no. 2, pp. 1146–1158, 2017.
- [16] M. Ji, A. M. Tulino, J. Llorca, and G. Caire, "On the average performance of caching and coded multicasting with random demands," in *IEEE International Symposium on Wireless Communications Systems (ISWCS)*. IEEE, 2014, pp. 922–926.
- [17] —, "Order-optimal rate of caching and coded multicasting with random demands," *IEEE Transactions on Information Theory*, vol. 63, no. 6, pp. 3923–3949, 2017.
- [18] J. Zhang, X. Lin, and X. Wang, "Coded caching under arbitrary popularity distributions," *IEEE Transactions on Information Theory*, vol. 64, no. 1, pp. 349–366, 2018.
- [19] S. Jin, Y. Cui, H. Liu, and G. Caire, "Structural properties of uncoded placement optimization for coded delivery," *arXiv preprint arXiv:1707.07146*, 2017.
- [20] A. M. Daniel and W. Yu, "Optimization of heterogeneous coded caching," *arXiv preprint arXiv:1708.04322*, 2017.
- [21] S. A. Saberali, L. Lampe, and I. Blake, "Full characterization of optimal uncoded placement for the structured clique cover delivery of nonuniform demands," *arXiv preprint arXiv:1804.00807*, 2018.
- [22] E. Ozfatura and D. Gunduz, "Uncoded caching and cross-level coded delivery for non-uniform file popularity," *arXiv preprint arXiv:1802.01135*, 2018.
- [23] H. Ding and L. Ong, "An improved caching scheme for nonuniform demands and its optimal allocation," in *EEE International Conference on Computer and Communications (ICCC)*. IEEE, 2017, pp. 389–393.
- [24] P. Quinton, S. Sahraei, and M. Gastpar, "A novel centralized strategy for coded caching with non-uniform demands," *arXiv preprint arXiv:1801.10563*, 2018.
- [25] K. Lee, M. Lam, R. Pedarsani, D. Papailiopoulos, and K. Ramchandran, "Speeding up distributed machine learning using codes," *IEEE Transactions on Information Theory*, vol. 64, no. 3, pp. 1514–1529, 2018.
- [26] S. Li, M. A. Maddah-Ali, and A. S. Avestimehr, "Coded mapreduce," in *2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, 2015, pp. 964–971.
- [27] S. El Rouayheb and K. Ramchandran, "Fractional repetition codes for repair in distributed storage systems," in *2010 48th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, 2010, pp. 1510–1517.

**Saeid Sahraei** received his Ph.D. and M.S. in 2018 and 2013, both from the School of Computer and Communication Sciences, École Polytechnique Fédérale de Lausanne (EPFL). Prior to that, he obtained his B.S. in Electrical Engineering from Sharif University of Technology in 2010. He is currently a postdoctoral scholar at the Department of Electrical Engineering, University of Southern California (USC). His research interests are in information theory, coding theory, combinatorial optimization and computational complexity.

Dr. Sahraei received the SNSF Early Postdoc.Mobility Fellowship in 2018 and the EPFL EDIC Doctoral Fellowship in 2013.

**Pierre Quinton** received a master's degree in Data Science from École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland, in 2019. During this time he worked on information theory and coding with Prof. Michael Gastpar and Prof. Emre Telatar. He is currently pursuing his PhD degree at EPFL under the supervision of Prof. Emre Telatar in the Information Theory Laboratory within the School of Computer and Communication Sciences. His main research interests lie within the areas of information theory, probability theory, and statistics.

**Michael Gastpar** (S'99–M'03–SM'14–F'17) received the Dipl. El.-Ing. degree from the Eidgenössische Technische Hochschule (ETH), Zürich, Switzerland, in 1997, the M.S. degree in electrical engineering from the University of Illinois at Urbana-Champaign, Urbana, IL, USA, in 1999, and the Doctorat ès Science degree from the Ecole Polytechnique Fédérale (EPFL), Lausanne, Switzerland, in 2002. He was also a student in engineering and philosophy at the Universities of Edinburgh and Lausanne.

During the years 2003–2011, he was an Assistant and tenured Associate Professor in the Department of Electrical Engineering and Computer Sciences at the University of California, Berkeley. Since 2011, he has been a Professor in the School of Computer and Communication Sciences, Ecole Polytechnique Fédérale (EPFL), Lausanne, Switzerland. He was also a professor at Delft University of Technology, The Netherlands, and a researcher with the Mathematics of Communications Department, Bell Labs, Lucent Technologies, Murray Hill, NJ. His research interests are in network information theory and related coding and signal processing techniques, with applications to sensor networks and neuroscience.

Dr. Gastpar received the IEEE Communications Society and Information Theory Society Joint Paper Award in 2013 and the EPFL Best Thesis Award in 2002. He was an Information Theory Society Distinguished Lecturer (2009–2011), an Associate Editor for Shannon Theory for the IEEE TRANSACTIONS ON INFORMATION THEORY (2008–2011), and he has served as Technical Program Committee Co-Chair for the 2010 International Symposium on Information Theory, Austin, TX.