

# A novel centralized coded caching scheme with coded prefetching

Jesús Gómez-Vilardebó, *Senior Member, IEEE*

**Abstract**—For the caching problem, when the number of files is no larger than that of users, the best known rate-memory region is achieved by memory sharing between the rate-memory pairs obtained by three schemes: the scheme proposed by Yu et al., the scheme proposed by Gomez-Vilardebo and the scheme proposed by Tian-Chen. While the first two schemes operate on the binary field, the Tian-Chen scheme makes use of a finite field of order  $2^m$  with, in some situations,  $m \geq K \log_2(N)$  for a caching systems with  $K$  users and  $N$  files. The practical implications of this increase in the size of the field are equivalent to an increase, by a factor of  $m$ , in the number of subfile partitions required. We propose a novel caching scheme that approaches the rate-memory region achieved by the Tian-Chen scheme as the number of users in the system increases, which only requires a field of order  $2^2$ .

**Index Terms**—Centralized coded caching, network coding, index coding.

## I. INTRODUCTION

Mobile data traffic has been increasing exponentially fast for the last 5-6 years. The advent of multimedia-capable devices such as smartphones and tablets at economical cost, as well as the success of content based applications such as YouTube, Facebook, Netflix, and network gaming, are the main reasons behind this growth. All signs indicate that this trend is likely to continue, for instance, to support HQ video on demand streaming applications. This constantly increasing demand of contents has been guiding the efforts of the industry and research communication communities. One of the most promising techniques proposed so far is “content caching” [1]. The most direct application of the content caching concept is to combat peak hour traffic in content delivery services. The simplest “uncoded” caching solutions [2], [3] work as follows: during low traffic periods, the storage resources available at the edge of the network are filled out with popular content. Then, whenever a user requests a content that is available at the edge cache, the content is served directly, thus, completely alleviating the backhaul network. This uncoded caching solution is only optimal in single cache systems. However, for multiple caches, the seminal work in [1] showed that important gains can be obtained by a new coded caching strategy. Specifically, in [1] authors show that, besides the local caching gain that is obtained by placing contents at user caches before they are requested, it is possible to obtain a global

caching gain by creating broadcast opportunities. This is, by carefully choosing the content caches at different users, and using network coding techniques it is possible to transform the initial unicast network, where every user is requesting a different file, into a broadcast network, where every user requests exactly the same “coded” file, obtaining the new global caching gain.

The fundamental caching scheme developed in [1] was later extended to more realistic situations. The decentralized setting was considered in [4], non-uniform demands in [5]–[7], on-line coded caching in [8], and hierarchical cache network were considered in [9], among others. In addition, new schemes pushing further the fundamental limits of caching systems have appeared in [10]–[18]. There have been also efforts to obtain theoretical lower bounds on the delivery rate. The cut-set bound was studied in [1]. A tighter lower bound was obtained in [19]. Through a computational approach a lower-bound for the special case with 3 files and 3 users is derived in [20]. Other lower bounds have appeared in [12], [21], [22].

Perhaps the most important challenge that needs further attention before coded caching solutions can be incorporated into real systems is what is known as the subpacketization level [23]–[25]. The theoretical gains of coded caching are demonstrated with strategies that require dividing each cached file into a number of parts that grows exponentially with the number of caches in the system. As an example, for the particular case where users store half of the full server data base, the centralized strategy in [1] which constitutes a fundamental building block for any of the current known extensions, needs to divide each packet into  $10^{14}$  subfiles with only 50 cache-users. The practical implications of such subpacketization levels are important in terms of the overall system complexity, signaling overhead and delay [26].

A particularly interesting situation is when there are more users  $K$  than files  $N$ . It was shown in [7] that a near optimal caching strategy consists in dividing the files into groups with similar popularity, and then applying the coded caching strategy to each group separately. Since the amount of users in each groups remains the same, when there are many groups, the cache size dedicated to each group is small as well as the number of files per user in each group.

If there are more users than files and cache memories are small, some of the best rate-memory pairs known are achieved by the Tian-Chen coding scheme [16]. This coding scheme, however, makes use, in some situations, of a finite field of at least order  $N^K$ . The practical implications of this increase in the size of the field are equivalent to an increase, by a factor of  $K \log_2 N$ , in the number of subfile partitions required compared to the scheme in [1]. We propose a novel caching

Manuscript received ; revised ; accepted; approved by IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS Editor . This work was supported in part by the Catalan Government under Grant SGR2017-1479, and in part by the Spanish Government under Grant TEC2013-44591-P (INTENSIV). The author is with the Centre Tecnològic de Telecomunicacions de Catalunya (CTTC/CERCA), 08860 Barcelona, Spain (e-mail: jesus.gomez@cttc.cat).

scheme that approaches the rate-memory region achieved by the Tian-Chen scheme as the number of users in the system increases, which only requires a finite field of order  $2^2$ . Moreover, instead of relying on the existence of some valid code as in [17], we provide an explicit combinatorial construction of the caching scheme, including, both, the prefetching, and delivery (broadcast and decoding) phases.

The rest of this paper is organized as follows. In Section II, we present the system model together with the more relevant previous results. In Section III, we summarize the main results. Section IV describes the caching scheme proposed. A detailed example is developed in Section V. Finally, conclusions are drawn in Section VI.

## II. SYSTEM MODEL AND PREVIOUS RESULTS

We consider a communication system with one server connected to  $K$  users, denoted as  $U_1, \dots, U_K$ , through a shared, error-free link. There is a database at the server with  $N$  files, each of length  $F$  bits, denoted as  $W_1, \dots, W_N$ . Each user is equipped with a local cache of capacity  $MF$  bits,  $M \in [0, N]$ , and is assumed to request only one full file. Here, we consider the special case where there are more users than files  $N \leq K$ .

We consider the communication model introduced in [1]. The caching system operates in two phases: the *placement phase* and the *delivery phase*. In the placement phase, users have access to the server database, and each user fills their cache. We allow coding in the prefetching phase. Thus, at user  $U_k$ , the caching function maps the database  $\{W_1, \dots, W_N\}$  to the cache content  $\mathcal{M}_k$ . Then, each user  $U_k$  requests a single full file  $W_{\mathbf{d}(k)}$ , where  $\mathbf{d} = (\mathbf{d}(1), \dots, \mathbf{d}(K))$  denotes the demand vector. We denote the number of distinct requests in  $\mathbf{d}$  as  $N_e(\mathbf{d})$ . In the delivery phase, only the server has access to the database. After being informed of the user demands, the server transmits a signal  $Y$  of size  $RF$  bits over the shared link to satisfy all user requests simultaneously. The signal  $Y$  is a function of the demand vector  $\mathbf{d}$ , all the files in the data base  $W_1, \dots, W_N$ , and the content in the user caches  $\mathcal{M} = \{\mathcal{M}_1, \dots, \mathcal{M}_k\}$ . Using the local cache content and the received signal  $Y$ , each user  $U_k$  reconstructs its requested file  $W_{\mathbf{d}(k)}$ .

Let  $\mathcal{D} = \{1, \dots, N\}^K$ , for a caching system  $(M, N, K)$ , given a particular prefetching  $\mathcal{M} = \{\mathcal{M}_1, \dots, \mathcal{M}_k\}$  and a particular demand  $\mathbf{d}$ , we say that communication rate  $R$  is achievable if and only if there exists a message  $Y$  of length  $RF$  bits such that every user  $U_k$  is able to reconstruct its desired file  $W_{\mathbf{d}(k)}$ . For a particular prefetching  $\mathcal{M}$  and demand  $\mathbf{d}$ , we denote the achievable rate as  $R(\mathbf{d}, \mathcal{M})$ . Then, the rate needed for the worst demand is given by  $R^*(\mathcal{M}) = \max_{\mathbf{d} \in \mathcal{D}} R(\mathbf{d}, \mathcal{M})$ . Finally, we define the rate-memory pair  $(R^*, M)$ , or the rate-memory trade-off function  $R^*(M)$  as the minim rate  $R^*$  for different memory constraints  $M$ , i.e. we aim to find

$$R^* = \min_{\mathcal{M}} R^*(\mathcal{M})$$

where the minimization is over all caching schemes  $\mathcal{M}$  satisfying the memory load constrain  $M$ . Observe, that if users have no caching capacity  $M = 0$ , the server needs to send the full requested files and thus, the worst demand rate is

$R^* = N$ . Instead, if users can have a complete copy of the server's database  $M = N$ , then no information needs to be transmitted from the server  $R^* = 0$ .

As far as we known, for the centralized caching setting, the best known explicitly characterization of the rate-memory trade-off can be obtained by memory sharing between three rate-memory pair families: the rate-memory pairs presented in [27, Corollary 1.1],

$$(R_{\text{GV}}^*, M_{\text{GV}}) = \left( N - \frac{N}{K} \frac{N+1}{q+1}, \frac{N}{Kq} \right) \quad (1)$$

for  $q \in \{1, \dots, N\}$ ; the rate-memory pairs obtained in [15, Corollary 1],

$$(R_{\text{YU}}^*, M_{\text{YU}}) = \left( \frac{\binom{K}{r+1} - \binom{K-\min\{K, N\}}{r+1}}{\binom{K}{r}}, \frac{rN}{K} \right) \quad (2)$$

for  $r \in \{0, 1, \dots, K\}$ , and the rate-memory pairs obtained in [16, Theorem 1]

$$(R_{\text{TC}}^*, M_{\text{TC}}) = \left( N \left( 1 - \frac{r}{K} \right), \frac{r}{K} + (N-1) \frac{r}{K} \frac{r-1}{K-1} \right), \quad (3)$$

with  $r \in \{0, 1, \dots, K\}$ .

The rate-memory pairs in (2) include as special case the rate-memory pair found in [10] for  $M = \frac{N}{K}$ . In addition, both (1) and (3) include as special case the rate-memory pair  $(R^*, M) = (N - \frac{N}{K}, \frac{1}{K})$  obtained in [14]. Recently, a method to obtain new rate-memory pairs has been developed in [17]. However, no explicit characterization of the achievable rate-memory pairs is provided. The optimal rate-memory trade-off for a caching systems remains an open problem. There have been efforts also to obtain theoretical lower bounds on the delivery rate. The cut-set bound was studied in [1]. A tighter lower bound was obtained in [19]. Through a computational approach a lower-bound for the special case  $N = K = 3$  is derived in [20]. Other lower bounds have appeared in [12], [21], [22].

Yu et al. [15] and Tian-Chen [16] coded caching schemes, both define exactly the same subfile partitioning and assignment strategies. Each file is divided into  $s = \binom{K}{r}$  subfiles, out of which to every user  $s_a = \binom{K-1}{r-1}$  subfiles are assigned in a way that each subfile is assigned to  $r$  users. Since there are  $N$  files, the total number of subfiles assigned to a user is  $Ns_a$ . Let  $s_{\bar{a}} = s - s_a$  denote the number subfiles of every file not assigned to a user. The prefetching and delivery strategies are however quite different for both schemes. In Yu et al. uncoded prefetching is proposed. At each user, the subfiles assigned are directly cached without further encoding. The number of subfiles cached is thus  $c_{\text{YU}} = s_a$ . The objective of the delivery strategy is to efficiently transmit the  $s_{\bar{a}}$  requested subfiles missing at each user. Observe that a "greedy" uncoded transmission strategy could be to simply broadcast the  $s_{\bar{a}}$  requested subfiles missing at each user. This "greedy" strategy requires a total of  $Ks_{\bar{a}}$  broadcasted subfiles. The Yu et al. broadcasting strategy finds multicast opportunities that allows a set of users to simultaneously recover their requested subfiles, as a result at most  $b_{\text{YU}} = \binom{K}{r+1} = \frac{K}{r+1} s_{\bar{a}}$  coded broadcasted subfiles are required. The Tian-Chen scheme, instead, make use of coded

prefetching. At each user, coded cached subfiles are obtained as linearly independent combinations, in some finite field, of the  $s_a$  subfiles assigned to the user. The total number of coded cached subfiles at each user is

$$c_{TC} = (N-1) \binom{K-2}{r-2} + \binom{K-1}{r-1}.$$

For any coded prefetching solution, the delivery strategy needs to, not only deliver the  $s_{\bar{a}}$  requested subfiles missing at each user, but also to allow decoding the  $s_a$  requested subfiles encoded in the cache. The Tian-Chen delivery scheme achieves these two objectives but, additionally, requires each user to completely decode the  $Ns_a$  subfiles encoded in its cache<sup>1</sup>. Observe that a “greedy” delivery strategy in this case, could be to broadcast to each user the  $s_{\bar{a}}$  missing requested subfiles uncoded, together with  $Ns_a - c_{TC}$  linearly independent combination of the subfiles coded in the cache, that together with the  $c_{TC}$  coded subfiles already in the cache could allow the decoding of every subfile in the cache. For this greedy scheme, the required total number of broadcasted subfiles is

$$b_{GC} = K(s_{\bar{a}} + Ns_a - c_{TC}) = K \left( r \frac{N-1}{K-1} + 1 \right) s_{\bar{a}}.$$

Instead, the Tian-Chen coded caching strategy is able to solve the same problem, broadcasting only  $b_{TC} = Ns_{\bar{a}}$  coded subfiles. To obtain this result, the authors in [16], [17] employed codes, in particular maximum distant separable (MDS) codes and/or rank metric codes, operating in a finite field of at least order

$$\min \left( N^K, 2 \binom{K-1}{r-1} N - \binom{K-2}{r-1} (N-1) \right). \quad (4)$$

This is in contrast to the combinatorial description of schemes presented in [15] and [27], which only require binary operations.

### III. MAIN RESULT

The following theorem presents the delivery rate obtained by the proposed caching scheme for a particular demand  $\mathbf{d}$ .

**Theorem 1.** *For a caching problem with  $K$  users,  $N$  files, and local cache size  $MF$  bits at each user. Given a particular demand  $\mathbf{d}$ , let  $N_e(\mathbf{d})$  be the number of distinct file requests, the proposed strategy achieves the rate-memory pairs*

$$R_{TC+} = \min(N_e(\mathbf{d}) + 1, N) \left( 1 - \frac{r}{K+1} \right)$$

$$M_{TC+} = \frac{r}{K+1} + (N-1) \frac{r}{K+1} \frac{r-1}{K}$$

for  $r \in \{0, \dots, K+1\}$ . Moreover, the new coded caching strategy operates in a finite field  $\mathbb{F}_4$ .

We prove this result in the following section by describing the new caching scheme. The delivery rate presented in Theorem 1 is valid for any  $K$  and  $N$ . However, it is particularly useful for  $K \geq N$ , as we detail in the next corollaries.

<sup>1</sup>This additional condition is probably limiting the performance of the strategy, however is a key ingredient in their result

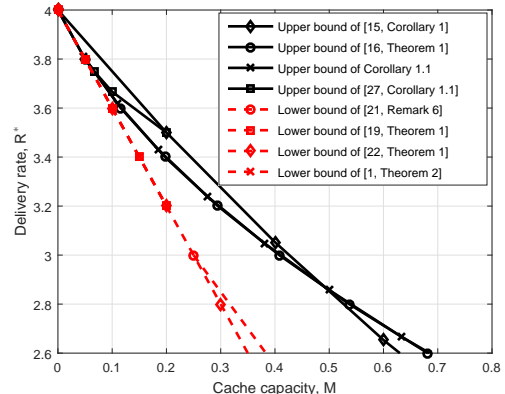


Figure 1: Rate-memory functions of the proposed scheme compared with existing schemes and lower bounds in the literature for  $N = 4$  and  $K = 20$ .

The next corollary provides the rate-memory region for the worst demand as a direct consequence of Theorem 1.

**Corollary 1.1.** *For a caching problem with  $K$  users and  $N$  files, local cache size of  $MF$  bits at each user, given that  $N_e(\mathbf{d}) \leq K$ , the proposed strategy achieves the worst demand rate-memory pairs  $(R^*, M) = (R_{TC+}^*, M_{TC+})$ , with*

$$R_{TC+}^* = \min(N, K+1) \left( 1 - \frac{r}{K+1} \right)$$

$$M_{TC+} = \frac{r}{K+1} + (N-1) \frac{r}{K+1} \frac{r-1}{K}$$

for  $r \in \{0, \dots, K+1\}$ .

**Remark 1.1.** *Observe that, if  $K \geq N$  the worst demand rate-memory function in Corollary 1.1 for  $K$  users coincides with the Tian-Chen rate-memory pair in (3) for  $K+1$  users.*

The following corollary compares both rate-memory functions. The proof is provided in the Appendix.

**Corollary 1.2.** *Let  $R_{TC+}^*(M)$  be the worst demand rate-memory function achieved by memory sharing between the rate-memory pairs in Corollary 1.1, and let  $R_{TC}^*(M)$  be the worst demand rate-memory function of the Tian-Chen scheme. For any number of files  $N$  and users  $K$ , satisfying  $K \geq N$  and  $M \in [0, N]$ , we have*

$$\frac{R_{TC+}^*(M)}{R_{TC}^*(M)} \leq 1 + \frac{1}{(K-1)(K+1)}.$$

Moreover, we have  $R_{TC+}^*(M) = R_{TC}^*(M)$  at

$$M = \frac{r+1}{K+1} \left( 1 + (N-1) \frac{r}{K} \right)$$

for all  $r \in \{0, \dots, K\}$ , as well as, at the intervals

$$0 \leq M \leq \frac{1}{K+1}, \frac{K-1}{K+1}N + \frac{1}{K+1} \leq M \leq N.$$

Observe that the proposed scheme matches Tian-Chen result for very low and very high caching ratios. For a system with  $N = 2$  files and  $K = 3$  users, Fig. 2 in the appendix, shows

the regimes where the proposed scheme and the Tian-Chen scheme are equal and those where they differ. In Fig. 1 for a caching system with  $N = 4$  files and  $K = 20$  users, we compare the rate-memory pairs in Corollary 1.1 to the best known rate-memory functions in the literature, i.e. the rate-memory pairs in [15, Corollary 1], [27, Corollary 1.1] and [16, Theorem 1]. In addition, we depict for comparison, the cut set lower bound [1, Theorem 2], the information-theoretical lower bound obtained in [19, Theorem 1], and the lower bounds recently appeared in [22, Theorem 1] and [21, Remark 6]. Observe that, as expected from the result in Corollary 1.2, the new proposed scheme obtains a rate-memory function overlapping almost exactly the rate-memory function in [16], [17]. Moreover, the strategy presented here only requires a finite field of order  $2^2$ . Instead, the direct evaluation of (4), for the rate-memory pair associated to  $r = 11$ , shows us that the Tian-Chen scheme requires a finite field of at least size  $10^5$ .

#### IV. PROPOSED CACHING SCHEME

We prove Theorem 1 in the following subsections. First, we introduce the main ideas of the new scheme with three examples. Then, we describe a caching scheme that achieves the Tian-Chen rate-memory pairs only if all files are requested by some user. This scheme fails to satisfy other type of demands where some files are not requested. Next, we show how the same scheme can satisfy any demand if configured for  $K+1$  caches instead of  $K$ . To help following every step of the description of the general scheme, we provide an extended example in the next section.

##### A. Motivation Examples

Here we present three examples. The first example considers the case when all files are requested but only by one user. The second example considers the case when not all files are requested. Finally, the third example considers the case when some files are requested by more than one user.

**Example 1:** Consider a caching system with  $K = 4$  users  $U_1, \dots, U_4$  and  $N = 4$  files A, B, C, D. Let us fix  $r = 3$ , then each file is partitioned into  $\binom{4}{3} = 4$  subfiles, and each subfile is assigned to  $r = 3$  users. To make explicit the users to which each subfile is assigned, we denote subfiles as  $A_{1,2,3}$ ,  $A_{1,2,4}$ ,  $A_{1,3,4}$ , and  $A_{2,3,4}$ . Then, the subfiles assigned to user  $U_1$  are

$$A_{1,2,3}, A_{1,2,4}, A_{1,3,4}, B_{1,2,3}, B_{1,2,4}, B_{1,3,4} \\ C_{1,2,3}, C_{1,2,4}, C_{1,3,4}, D_{1,2,3}, D_{1,2,4}, D_{1,3,4}$$

For the delivery strategy, we consider the uncoded transmission of the requested subfiles not encoded in the cache. For instance, if users  $U_1, U_2, U_3, U_4$  request A, B, C, D, respectively. Then, the server broadcasts

$$A_{2,3,4}, B_{1,3,4}, C_{1,2,4}, D_{1,2,3}.$$

Observe that, in general, there are  $Ks_{\bar{a}}$  of these subfiles. Observe also that, since  $K = N$ , this coincides with the total number of broadcasted subfiles for the Tian-Chen schemes, i.e.  $b_{TC} = Ns_{\bar{a}}$ . To design the prefetching scheme, we impose, as in the Tian-Chen scheme, the decoding of all

the subfiles in the cache for any possible demand. For user  $U_1$ , observe that, for each of the files B, C, D not requested by user  $U_1$ , hereafter interfering files, the server broadcasts one subfile coded in the cache of user  $U_1$ , i.e.  $B_{1,3,4}, C_{1,2,4}, D_{1,2,3}$ . Taking this into account, we design a set of coded subfiles only XORing subfiles of one file, such that after receiving any subfile of this file, the rest can be decoded. For this particular examples, for file B, we could fill the cache at user  $U_1$  with

$$B_{1,2,3} \oplus B_{1,2,4}, B_{1,2,3} \oplus B_{1,3,4}.$$

Observe that after receiving  $B_{1,3,4}$ , user  $U_1$  can obtain all the subfiles of file B coded in his cache. For symmetry, the equivalent coded subfiles are cached for other other three files.

We still need to obtain the requested subfiles encoded in the cache. However, given that all interfering subfiles are already known, any set of coded subfiles, each XORing together subfiles of all  $N$  files, and satisfying that each subfile is XORed only once, will work. For instance, we can use

$$A_{1,2,3} \oplus B_{1,2,3} \oplus C_{1,2,3} \oplus D_{1,2,3} \\ A_{1,2,4} \oplus B_{1,2,4} \oplus C_{1,2,4} \oplus D_{1,2,4} \\ A_{1,3,4} \oplus B_{1,3,4} \oplus C_{1,3,4} \oplus D_{1,3,4}$$

Observe that we only needed binary operations.

**Example 2:** Next, yet assuming that each files is requested by one user at most, we address the situation when not all files are requested. Consider a caching system with  $K = 3$  and  $N = 4$  files. Observe that if we assume that there is a 4th “virtual” user that always requests the file not requested by the others, then we can use the coded caching strategy presented in Example 1 for  $K = 4$  users. In general, for a system with  $K$  users, we can achieve the Tian-Chen rate-memory curve for  $K+1$  users.

**Example 3:** Next, we address the situation where some files are requested by more than one user. For this case, the binary operation field would not be sufficient. We focus on a caching system with  $N = 2$  files A, B,  $K = 3$  users, and  $r = 2$ . By applying the coded prefetching scheme introduced in Example 1, we obtain

$U_1$	$U_2$	$U_3$
$A_{1,2} \oplus B_{1,2}$	$A_{1,2} \oplus B_{1,2}$	$A_{1,3} \oplus B_{1,3}$
$A_{1,3} \oplus B_{1,3}$	$A_{2,3} \oplus B_{2,3}$	$A_{2,3} \oplus B_{2,3}$
$A_{1,2} \oplus A_{1,3}$	$A_{1,2} \oplus A_{2,3}$	$A_{1,3} \oplus A_{2,3}$
$B_{1,2} \oplus B_{1,3}$	$B_{1,2} \oplus B_{2,3}$	$B_{1,3} \oplus B_{2,3}$

Suppose file A is requested by users  $U_1$  and  $U_3$ , whereas file B is requested by user  $U_2$ . Then, the uncoded broadcasting strategy used for Example 1 transmits all the requested subfiles not encoded at the user requesting them, i.e.  $A_{1,2}, A_{2,3}, B_{1,3}$ . Notice that these subfiles are sufficient for every user to obtain the requested subfiles and decode all the subfiles in the cache. However, since there are more users than files  $K > N$ , the uncoded broadcasted scheme needs  $Ks_{\bar{a}}$  subfiles transmission instead of the  $Ns_{\bar{a}}$  required by the Tian-Chen scheme. To see the limitations of the uncoded broadcasted scheme, observe that upon recovering all the subfiles in the cache user  $U_1$  has  $A_{1,2}$  and user  $U_2$  has  $A_{2,3}$  and thus, broadcasting  $A_{1,2} \oplus A_{2,3}$  would be sufficient for both to recover the missing subfiles. However, the coded subfile  $A_{1,2} \oplus A_{2,3}$

is already at user  $U_2$  cache and is, thus, useless to him. To resolve this issue, we further divide each subfiles into two half-subfiles,  $A_{1,2} = \langle A_{1,2}^I, A_{1,2}^Q \rangle$  and consider the previous prefetching scheme over the half-subfiles. Then, the cache of users  $U_1$  and  $U_2$  read

$U_1$		$U_2$	
$A_{1,2}^I \oplus B_{1,2}^I$	$A_{1,2}^Q \oplus B_{1,2}^Q$	$A_{1,2}^I \oplus B_{1,2}^I$	$A_{1,2}^Q \oplus B_{1,2}^Q$
$A_{1,3}^I \oplus B_{1,3}^I$	$A_{1,3}^Q \oplus B_{1,3}^Q$	$A_{2,3}^I \oplus B_{2,3}^I$	$A_{2,3}^Q \oplus B_{2,3}^Q$
$A_{1,2}^I \oplus A_{1,3}^I$	$A_{1,2}^Q \oplus A_{1,3}^Q$	$A_{1,2}^I \oplus A_{2,3}^I$	$A_{1,2}^Q \oplus A_{2,3}^Q$
$B_{1,2}^I \oplus B_{1,3}^I$	$B_{1,2}^Q \oplus B_{1,3}^Q$	$B_{1,2}^I \oplus B_{2,3}^I$	$B_{1,2}^Q \oplus B_{2,3}^Q$

Now, observe that broadcasting the half-subfiles  $A_{1,2}^I \oplus A_{2,3}^I \oplus A_{2,3}^Q$ ,  $A_{1,2}^I \oplus A_{1,2}^Q \oplus A_{2,3}^Q$ ,  $B_{1,3}^I$ , and  $B_{1,3}^Q$ , allows every user to recover their requested subfiles. Observe that  $U_1$  and  $U_3$  upon receiving  $B_{1,3}^I$ , and  $B_{1,3}^Q$  can cancel all the interference from file B in their cache. Consequently,  $U_1$  also obtains  $A_{1,2}^I$  and  $A_{1,2}^Q$ , which finally allows him to obtain first  $A_{2,3}^Q$  from  $A_{1,2}^I \oplus A_{1,2}^Q \oplus A_{2,3}^Q$  and then  $A_{2,3}^I$  from  $A_{1,2}^I \oplus A_{2,3}^I \oplus A_{2,3}^Q$ .  $U_3$  proceeds similarly. Observe that  $U_2$  can also cancel the interference of file A using the broadcasted half-subfiles and the half-subfiles in their cache, by first computing

$$\begin{aligned} [A_{1,2}^I \oplus A_{2,3}^I] \oplus [A_{1,2}^I \oplus A_{2,3}^I \oplus A_{2,3}^Q] &= A_{2,3}^Q \\ [A_{1,2}^Q \oplus A_{2,3}^Q] \oplus [A_{1,2}^I \oplus A_{1,2}^Q \oplus A_{2,3}^Q] &= A_{1,2}^I \end{aligned}$$

and then

$$\begin{aligned} [A_{1,2}^Q \oplus A_{2,3}^Q] \oplus A_{2,3}^Q &= A_{1,2}^Q \\ [A_{1,2}^I \oplus A_{2,3}^I] \oplus A_{1,2}^I &= A_{2,3}^I. \end{aligned}$$

Finally, after removing all interference from file A,  $U_2$  obtains the requested subfiles  $B_{1,2}^I, B_{1,2}^Q, B_{2,3}^I, B_{2,3}^Q$  coded in its cache.

### B. Coded caching scheme if all files are requested

Let us first consider a caching system that can only serve demand vectors that include all files in the library.

1) *Prefetching scheme*: Let us define the set of users indexes as  $\mathcal{K} = \{1, \dots, K\}$  and the set of file indexes as  $\mathcal{F} = \{1, \dots, N\}$ . We partition each file  $W_f$ ,  $f \in \mathcal{F}$  into  $\binom{K}{r}$  non-overlapping subfiles of equal size<sup>2</sup>,  $W_{f,S}$ , one for each subset  $\mathcal{S}$  of  $r \in \{0, \dots, K\}$  users, i.e.  $\mathcal{S} \in \mathcal{T}(r)$  with

$$\mathcal{T}(r) = \{\mathcal{S} \subseteq \mathcal{K} : |\mathcal{S}| = r\}.$$

Let  $\mathcal{T}_k(r) = \{\mathcal{S} \in \mathcal{T}(r) : k \in \mathcal{S}\}$  be the set of subsets of  $r$  users that include a particular user  $U_k$ . User  $U_k$  caches the binary sum of the  $N$  subfiles associated to the same subset  $\mathcal{S}$

$$Z_{\mathcal{S}} = \bigoplus_{f \in \mathcal{F}} W_{f,S} \quad (5)$$

for all  $\mathcal{S} \in \mathcal{T}_k(r)$ . Consequently, each user caches  $\binom{K-1}{r-1}$  coded subfiles. If  $r = 1$  these are all the coded cached subfiles need at the caches. If  $r \geq 2$ , we need more coded cached

<sup>2</sup>We define  $\binom{n}{k} = 0$  if  $n < k$ .

subfiles. Specifically, we need user  $U_k$  to be able to obtain from its cache the coded subfiles

$$Z_{f,S^-} = \bigoplus_{s \in \mathcal{K} \setminus S^-} W_{f,S^- \cup s} \quad (6)$$

for every file  $f \in \mathcal{F}$ , and every subset  $S^- \in \mathcal{T}_k(r-1)$ . As we show next, user  $U_k$  does not need to cache them all. Specifically, user  $U_k$  arbitrarily selects one file index  $g$ , and one user index  $l$ , and caches only the coded subfiles  $Z_{f,S^-}$  in (6) for all files  $f \in \mathcal{F} \setminus g$  and all sets  $S^- \in \mathcal{T}_k(r-1)$  satisfying  $l \notin S^-$ . The total number of coded subfiles  $Z_{f,S^-}$  cached at user  $U_k$  is then  $(N-1) \binom{K-2}{r-2}$ .

Then, whenever user  $U_k$  needs  $Z_{f,S^-}$  for any file  $f \in \mathcal{F} \setminus g$  and any subset  $S^- \in \mathcal{T}_k(r-1)$  satisfying  $l \in S^-$ , he only needs to XOR the coded subfiles  $Z_{f, \{S^- \setminus l\} \cup v}$  for all  $v \in \mathcal{K} \setminus S^-$ . Observe that since  $l \notin \{S^- \setminus l\} \cup v$ , this coded subfiles are cached at user  $U_k$ . Then,  $Z_{f,S^-}$

$$\begin{aligned} &= \bigoplus_{v \in \mathcal{K} \setminus S^-} Z_{f, \{S^- \setminus l\} \cup v} \\ &= \bigoplus_{v \in \mathcal{K} \setminus S^-} \bigoplus_{s \in \mathcal{K} \setminus \{\{S^- \setminus l\} \cup v\}} W_{f, \{S^- \setminus l\} \cup v \cup s} \\ &= \bigoplus_{v \in \mathcal{K} \setminus S^-} W_{f, S^- \cup v} \oplus \bigoplus_{v \in \mathcal{K} \setminus S^-} \bigoplus_{s \in \mathcal{K} \setminus \{S^- \cup v\}} W_{f, \{S^- \setminus l\} \cup v \cup s} \\ &= \bigoplus_{v \in \mathcal{K} \setminus S^-} W_{f, S^- \cup v} \end{aligned} \quad (7)$$

where equality (7) follows since in the previous equality  $W_{f, \{S^- \setminus l\} \cup v \cup s}$  appear twice, and thus cancel out.

To obtain  $Z_{g,S^-}$  for any set  $S^- \in \mathcal{T}_k(r-1)$ , user  $U_k$  XORs the coded cached subfiles  $Z_{S^- \cup s}$  for all  $\mathcal{K} \setminus S^-$  and the coded cached subfiles  $Z_{f,S^-}$  for all  $f \in \mathcal{F} \setminus g$ , as follows

$$\begin{aligned} Z_{g,S^-} &= \bigoplus_{s \in \mathcal{K} \setminus S^-} Z_{S^- \cup s} \oplus \bigoplus_{f \in \mathcal{F} \setminus g} Z_{f,S^-} \\ &= \bigoplus_{s \in \mathcal{K} \setminus S^-} Z_{S^- \cup s} \oplus \bigoplus_{f \in \mathcal{F}} Z_{f,S^-} \oplus Z_{g,S^-} \\ &= \bigoplus_{s \in \mathcal{K} \setminus S^-} Z_{S^- \cup s} \oplus \bigoplus_{s \in \mathcal{K} \setminus S^-} Z_{S^- \cup s} \oplus Z_{g,S^-} \end{aligned} \quad (8)$$

where (8) follows, since

$$\begin{aligned} \bigoplus_{f \in \mathcal{F}} Z_{f,S^-} &= \bigoplus_{f \in \mathcal{F}} \bigoplus_{s \in \mathcal{K} \setminus S^-} W_{f,S^- \cup s} \\ &= \bigoplus_{s \in \mathcal{K} \setminus S^-} \bigoplus_{f \in \mathcal{F}} W_{f,S^- \cup s} \\ &= \bigoplus_{s \in \mathcal{K} \setminus S^-} Z_{S^- \cup s}. \end{aligned}$$

Finally, because each subfile has  $F/\binom{K}{r}$  bits, the required cache load at users equals  $MF$ , with

$$M = \frac{\binom{K-1}{r-1} + (N-1) \binom{K-2}{r-2}}{\binom{K}{r}} = \frac{r}{K} + (N-1) \frac{r}{K} \frac{r-1}{K-1}.$$

The server further divides each subfile  $W_{f,S}$  into two half-subfiles of equal size  $W_{f,S}^I$  and  $W_{f,S}^Q$ . This can be seen, equivalently, as extending the field of the system to  $\mathbb{F}_{2^2}$ . Let us write

a subfile from its half-subfiles as  $W_{f,S} = \langle W_{f,S}^I, W_{f,S}^Q \rangle$ . At the caches, users use the same partition function to also divide the coded cached subfiles into two coded half-subfiles  $Z_S = \langle Z_S^I, Z_S^Q \rangle$  and  $Z_{f,S^-} = \langle Z_{f,S^-}^I, Z_{f,S^-}^Q \rangle$ . Then, the coded half-subfiles are given by

$$\begin{aligned} Z_S^I &= \bigoplus_{f \in \mathcal{F}} W_{f,S}^I, & Z_S^Q &= \bigoplus_{f \in \mathcal{F}} W_{f,S}^Q, \\ Z_{f,S^-}^I &= \bigoplus_{s \notin \mathcal{S}^-} W_{f,S^-}^I, & Z_{f,S^-}^Q &= \bigoplus_{s \notin \mathcal{S}^-} W_{f,S^-}^Q. \end{aligned}$$

2) Broadcasting scheme: Let us denote the binary sum of the two half-subfiles of subfile  $W_{f,S}$  as  $\bar{W}_{f,S} = W_{f,S}^I \oplus W_{f,S}^Q$  and the set of users requesting file  $W_f$  as  $\mathcal{K}(f)$ . The server, first, arbitrarily selects one user leader  $u_f \in \mathcal{K}(f)$  for each file  $f \in \mathcal{F}$ , and then broadcasts the coded subfile  $Y_{f,S} = \langle Y_{f,S}^I, Y_{f,S}^Q \rangle$  with half-subfiles, given by

$$Y_{f,S}^I = W_{f,S}^I \oplus \bigoplus_{s \in \mathcal{S} \cap \mathcal{K}(f)} \bar{W}_{f,u_f \cup \{s\}}, \quad (9)$$

$$Y_{f,S}^Q = \bar{W}_{f,S} \oplus \bigoplus_{s \in \mathcal{S} \cap \mathcal{K}(f)} W_{f,u_f \cup \{s\}}^Q \quad (10)$$

for every file  $f \in \mathcal{F}$  and every set  $\mathcal{S}$  of  $r$  users, i.e.  $\mathcal{S} \in \mathcal{T}(r)$ , excluding the user leader  $u_f$ , i.e.  $u_f \notin \mathcal{S}$ .

Because, there are  $N$  files, and  $\binom{K-1}{r}$  sets  $\mathcal{S} \in \{\mathcal{S} \in \mathcal{T}(r) : u_f \notin \mathcal{S}\}$  for each file, we need  $N \binom{K-1}{r}$  subfile length broadcast transmissions, thus, requiring the rate

$$R = \frac{N \binom{K-1}{r}}{\binom{K}{r}} = N \left(1 - \frac{r}{K}\right).$$

In order to facilitate the decoding process described next, we divide the broadcast of these subfiles into  $r+1$  consecutive phases  $p \in \{0, \dots, r\}$ . At phase  $p$ , for every file  $W_f$ , we broadcast the subfiles  $Y_{f,S}$  satisfying  $|\mathcal{S} \cap \mathcal{K}(f)| = p$ .

Recall the definition of  $\mathcal{T}_k(r)$  as the set of subsets of  $r$  users that include a particular user  $k$  and define its complement in  $\mathcal{T}(r)$  as  $\bar{\mathcal{T}}_k(r) = \{\mathcal{S} \in \mathcal{T}(r) : k \notin \mathcal{S}\}$ .

The decoding algorithm is summarized in Algorithm 1. User  $U_k$  divides the decoding of all the requested subfiles  $W_{d(k),S}$  for all  $\mathcal{S} \in \mathcal{T}(r)$  into two steps. First, user  $U_k$  obtains all subfiles  $W_{d(k),S}$  for all sets  $\mathcal{S} \in \mathcal{T}_k(r)$ . This subfiles are coded at some coded subfile in the cache of user  $U_k$ . Then, user  $U_k$  obtains the subfiles cached by other users, i.e. subfiles  $W_{d(k),S}$  for all sets  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$ .

3) *Decoding of requested subfiles coded in the cache of user  $U_k$* : We first describe how user  $U_k$  obtains subfiles  $W_{d(k),S}$  for all sets  $\mathcal{S} \in \mathcal{T}_k(r)$ . Recall that for any set  $\mathcal{S} \in \mathcal{T}_k(r)$ , user  $k$  caches

$$Z_S = W_{d(k),S} \oplus \bigoplus_{f \in \mathcal{F} \setminus \mathbf{d}(k)} W_{f,S}$$

and thus, can obtain  $W_{d(k),S}$  for all sets  $\mathcal{S} \in \mathcal{T}_k(r)$  by first obtaining all the interfering subfiles  $W_{f,S}$  for all  $f \in \mathcal{F} \setminus \mathbf{d}(k)$ , and then computing

$$W_{d(k),S} = Z_S \oplus \bigoplus_{f \in \mathcal{F} \setminus \mathbf{d}(k)} W_{f,S}. \quad (11)$$

The decoding of these interfering subfiles is divided into  $r$  consecutive phases  $p = 0, \dots, r-1$ . At phase  $p$ , for

---

**Algorithm 1** Decoding requested subfiles.

---

- 1) **Decode  $W_{f,S}$  for all  $\mathcal{S} \in \mathcal{T}_k(r)$  and  $f \in \mathcal{F} \setminus \mathbf{d}(k)$** 
    - a) **Input:**
      - $Z_S$  for all  $\mathcal{S} \in \mathcal{T}_k(r)$ ,
      - $Z_{f,S^-}$  for all  $f \in \mathcal{F}$  and  $\mathcal{S}^- \in \mathcal{T}_k(r-1)$ ,
      - $Y_{f,S}$  for all  $f \in \mathcal{F}$  and  $\mathcal{S} \in \mathcal{T}_k(r)$  satisfying  $u_f \notin \mathcal{S}$ .
    - b) **For all  $f \in \mathcal{F} \setminus \mathbf{d}(k)$ , and for  $p = 0, \dots, r-1$ ,**
      - i) **For all  $\mathcal{S} \in \mathcal{T}_k(r)$  with  $|\mathcal{S} \cap \mathcal{K}(f)| = p$  and  $u_f \in \mathcal{S}$** 
        - A) Compute  $C_{f,S \setminus u}$  for all  $u \in \mathcal{S} \cap \mathcal{K}(f)$  as (15)
        - B) Compute  $\Upsilon_{f,S \setminus u_f}$  as in (16)
        - C) Compute  $\Upsilon_{f,S}^I$  as in (17) and  $\Upsilon_{f,S}^Q$  as in (18)
        - D) Compute  $W_{f,S}$  as in (19)
      - ii) **For all  $\mathcal{S} \in \mathcal{T}_k(r)$  with  $|\mathcal{S} \cap \mathcal{K}(f)| = p$  and  $u_f \notin \mathcal{S}$** 
        - A) Compute  $W_{f,S}$  according to (13) and (14)
  - 2) **Decode  $W_{d(k),S}$  for all  $\mathcal{S} \in \mathcal{T}_k(r)$** 
    - a) **Input:**
      - $W_{f,S}$  for all  $\mathcal{S} \in \mathcal{T}_k(r)$  and  $f \in \mathcal{F} \setminus \mathbf{d}(k)$
      - $Z_S$  for all  $\mathcal{S} \in \mathcal{T}_k(r)$
    - b) **For all  $\mathcal{S} \in \mathcal{T}_k(r)$ , compute  $W_{d(k),S}$  as (11)**
  - 3) **Decode  $W_{d(k),S}$  for all  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$** 
    - a) **Input:**
      - $W_{d(k),S}$  for all  $\mathcal{S} \in \mathcal{T}_k(r)$ ,
      - $Y_{d(k),S}$  for all  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$  satisfying  $u_f \notin \mathcal{S}$
    - b) **Set  $f = \mathbf{d}(k)$**
    - c) **If  $k = u_f$ ,**
      - i) **For all  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$  compute  $W_{f,S}$  according to (25) and (26)**
    - d) **If  $k \neq u_f$ ,**
      - i) **For all  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$  satisfying  $u_f \in \mathcal{S}$  compute  $W_{f,S}$  according to (27) and (28)**
      - ii) **For all  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$  satisfying  $u_f \notin \mathcal{S}$  compute  $W_{f,S}$  according to (25) and (26)**
- 

every file  $f \in \mathcal{F} \setminus \mathbf{d}(k)$ , user  $U_k$  obtains subfile  $W_{f,S}$  for all subsets  $\mathcal{S} \in \mathcal{T}_k(r)$  with  $p$  users requesting file  $W_f$ , i.e. satisfying  $|\mathcal{S} \cap \mathcal{K}(f)| = p$ . Let us denote this family of sets as  $\mathcal{T}_{k,f}(r,p) = \{\mathcal{S} \in \mathcal{T}_k(r) : |\mathcal{S} \cap \mathcal{K}(f)| = p\}$ .

Consider first the decoding of the interfering subfiles at phase  $p = 0$ . This is, subfiles  $W_{f,S}$  for any set  $\mathcal{S} \in \mathcal{T}_{k,f}(r,0)$ . Observe that,  $u_f \notin \mathcal{S}$ , and thus, there is a broadcasted subfile  $Y_{f,S}$  associated to each of the sets  $\mathcal{S} \in \mathcal{T}_{k,f}(r,0)$ , with half-subfiles given by  $Y_{f,S}^I = W_{f,S}^I$ , and  $Y_{f,S}^Q = \bar{W}_{f,S}$ . Consequently, we can obtain  $W_{f,S} = \langle W_{f,S}^I, W_{f,S}^Q \rangle$  as

$$W_{f,S}^I = Y_{f,S}^I, \quad W_{f,S}^Q = Y_{f,S}^Q \oplus Y_{f,S}^I. \quad (12)$$

Next, consider the decoding of interfering subfiles at phases  $p, r > p \geq 1$ . This is, interfering subfiles  $W_{f,S}$  for any set  $\mathcal{S} \in \mathcal{T}_{k,f}(r,p)$  for  $p = 1, \dots, r-1$ . At phase  $p$ , user  $U_k$  decodes first all the interfering subfiles  $W_{f,S}$  also cached by the user leader, i.e.  $u_f \in \mathcal{S}$ . Then, the interfering subfiles  $W_{f,S}$  for all sets  $\mathcal{S} \in \mathcal{T}_{k,f}(r,p)$  not cached at the user leader  $u_f \notin \mathcal{S}$  can

be obtained as,  $W_{f,S} = \langle W_{f,S}^I, W_{f,S}^Q \rangle$  with

$$W_{f,S}^I = Y_{f,S}^I \oplus \bigoplus_{s \in \mathcal{S} \cap \mathcal{K}(f)} \bar{W}_{f,u_f \cup \{s\}}, \quad (13)$$

$$W_{f,S}^Q = Y_{f,S}^Q \oplus \bigoplus_{s \in \mathcal{S} \cap \mathcal{K}(f)} W_{f,u_f \cup \{s\}}^Q \oplus W_{f,S}^I. \quad (14)$$

Observe that, given that  $u_f \notin \mathcal{S}$ , the broadcasted subfile  $Y_{f,S}$  exists. Moreover, for all  $u \in \mathcal{S} \cap \mathcal{K}(f)$  the set  $u_f \cup \mathcal{S} \setminus u$  contains  $u_f$  and satisfies  $|u_f \cup \mathcal{S} \setminus u| = p$  and thus, user  $U_k$  has already obtained  $W_{f,u_f \cup \mathcal{S} \setminus u}$ .

Finally, we detail the decoding of interfering subfiles  $W_{f,S}$  for any set  $\mathcal{S} \in \mathcal{T}_{k,f}(r,p)$  with  $u_f \in \mathcal{S}$ . We summarize the decoding process first, and prove it later. Observe that for all  $u \in \mathcal{S} \cap \mathcal{K}(f)$ , we have  $|\{\mathcal{S} \setminus u\} \cap \mathcal{K}(f)| = p - 1$ , and thus from decoding phase  $p - 1$ , user  $U_k$  already has  $W_{f,\{\mathcal{S} \setminus u\} \cup v}$  for all  $v \in \mathcal{K}(f) \setminus \{\mathcal{S} \setminus u\}$ . This allows user  $U_k$  to compute

$$C_{f,S \setminus u} = Z_{f,S \setminus u} \oplus \bigoplus_{v \in \mathcal{K}(f) \setminus \{\mathcal{S} \setminus u\}} W_{f,\{\mathcal{S} \setminus u\} \cup v} \quad (15)$$

for all  $u \in \mathcal{S} \cap \mathcal{K}(f)$ . Next, observe that since  $u_f \in \mathcal{S}$ , and  $u_f \notin \mathcal{K}(f) \setminus \mathcal{S}$ , there is a broadcasted subfile  $Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}$  for all  $s \in \mathcal{K}(f) \setminus \mathcal{S}$ . This allows user  $U_k$  to compute

$$\Upsilon_{f,S \setminus u_f} = \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}. \quad (16)$$

Next, we partition  $\Upsilon_{f,S \setminus u_f}$  and  $C_{f,S \setminus u}$  for all  $u \in \mathcal{S} \cap \mathcal{K}(f)$  into half-subfiles, i.e.  $\Upsilon_{f,S \setminus u_f} = \langle \Upsilon_{f,S \setminus u_f}^I, \Upsilon_{f,S \setminus u_f}^Q \rangle$ ,  $C_{f,S \setminus u} = \langle C_{f,S \setminus u}^I, C_{f,S \setminus u}^Q \rangle$  and define  $\tilde{\Upsilon}_{f,S \setminus u_f} = \Upsilon_{f,S \setminus u_f}^I \oplus \Upsilon_{f,S \setminus u_f}^Q$ , and  $\tilde{C}_{f,S \setminus u} = C_{f,S \setminus u}^I \oplus C_{f,S \setminus u}^Q$ . Then, user  $U_k$  computes

$$\Upsilon_{f,S}^I = \Upsilon_{f,S \setminus u_f}^I \oplus C_{f,S \setminus u_f}^I \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} \tilde{C}_{f,S \setminus u} \quad (17)$$

$$\Upsilon_{f,S}^Q = \Upsilon_{f,S \setminus u_f}^Q \oplus \tilde{C}_{f,S \setminus u_f} \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} C_{f,S \setminus u}^Q \quad (18)$$

and obtains subfile  $W_{f,S} = \langle W_{f,S}^I, W_{f,S}^Q \rangle$  as

$$W_{f,S}^I = \begin{cases} \Upsilon_{f,S}^Q & \text{if } |\mathcal{K}(f) \setminus u_f| \text{ is odd} \\ \Upsilon_{f,S}^I & \text{otherwise} \end{cases} \quad (19)$$

$$W_{f,S}^Q = \begin{cases} \Upsilon_{f,S}^I & \text{if } |\mathcal{K}(f) \setminus u_f| \text{ is odd} \\ \Upsilon_{f,S}^I \oplus \Upsilon_{f,S}^Q & \text{otherwise} \end{cases}$$

In the following, we prove the correctness of this decoding process. First, observe that for all  $u \in \mathcal{S} \cap \mathcal{K}(f)$ , we have

$$\begin{aligned} C_{f,S \setminus u} &= Z_{f,S \setminus u} \oplus \bigoplus_{v \in \mathcal{K}(f) \setminus \{\mathcal{S} \setminus u\}} W_{f,\{\mathcal{S} \setminus u\} \cup v} \\ &= \bigoplus_{s \in \mathcal{K}(f) \setminus \{\mathcal{S} \setminus u\}} W_{f,\{\mathcal{S} \setminus u\} \cup s} \oplus \bigoplus_{v \in \mathcal{K}(f) \setminus \{\mathcal{S} \setminus u\}} W_{f,\{\mathcal{S} \setminus u\} \cup v} \\ &= \bigoplus_{s \in \mathcal{K}(f) \setminus \{\mathcal{S} \setminus u\}} W_{f,\{\mathcal{S} \setminus u\} \cup s} \\ &= W_{f,S} \oplus \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} W_{f,\{\mathcal{S} \setminus u\} \cup s} \end{aligned}$$

from which we can write

$$\bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} W_{f,\{\mathcal{S} \setminus u\} \cup s} = C_{f,S \setminus u} \oplus W_{f,S}. \quad (20)$$

Next, observe that for any  $s \in \mathcal{K}(f) \setminus \mathcal{S}$ , we can write  $Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}^I$

$$\begin{aligned} &= W_{f,\{\mathcal{S} \setminus u_f\} \cup s}^I \oplus \bigoplus_{u \in \{\{\mathcal{S} \setminus u_f\} \cup s\} \cap \mathcal{K}(f)} \bar{W}_{f,u_f \cup \{\{\mathcal{S} \setminus u_f\} \cup s\} \cup u} \\ &= W_{f,\{\mathcal{S} \setminus u_f\} \cup s}^I \oplus \bar{W}_{f,S} \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} \bar{W}_{f,\{\mathcal{S} \setminus u\} \cup s} \quad (21) \end{aligned}$$

where (21) follows since for  $u = s$ , we have  $u_f \cup \{\{\mathcal{S} \setminus u_f\} \cup s\} \setminus u = \mathcal{S}$ , and for  $u \neq s$ , since  $u_f \in \mathcal{S}$ , we have  $u_f \cup \{\{\mathcal{S} \setminus u_f\} \cup s\} \setminus u = \{\mathcal{S} \setminus u\} \cup s$ . Similarly, we can write  $Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}^Q$

$$= \bar{W}_{f,\{\mathcal{S} \setminus u_f\} \cup s} \oplus W_{f,S}^Q \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} W_{f,\{\mathcal{S} \setminus u\} \cup s}^Q. \quad (22)$$

Now, let  $\Upsilon_{f,S \setminus u_f}^I \triangleq \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}^I$ , and  $\Upsilon_{f,S \setminus u_f}^Q \triangleq \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}^Q$  then, it follows from (20) and (21) that

$$\begin{aligned} \Upsilon_{f,S \setminus u_f}^I &= C_{f,S \setminus u_f}^I \oplus W_{f,S}^I \oplus \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} \bar{W}_{f,S} \\ &\quad \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} \bar{C}_{f,S \setminus u} \oplus \bar{W}_{f,S} \quad (23) \end{aligned}$$

and, similarly, from (20) and (22) that

$$\begin{aligned} \Upsilon_{f,S \setminus u_f}^Q &= \bar{C}_{f,S \setminus u_f} \oplus \bar{W}_{f,S} \oplus \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} W_{f,S}^Q \\ &\quad \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} C_{f,S \setminus u}^Q \oplus W_{f,S}^Q. \quad (24) \end{aligned}$$

Next observe that

$$\begin{aligned} &\bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} W_{f,S} \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} W_{f,S} \\ &= \bigoplus_{s \in \mathcal{K}(f) \setminus u_f} W_{f,S} \end{aligned}$$

$$= \begin{cases} W_{f,S} & \text{if } |\mathcal{K}(f) \setminus u_f| \text{ is odd} \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

and thus  $\Upsilon_{f,S \setminus u_f}^I \oplus C_{f,S \setminus u_f}^I \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} \bar{C}_{f,S \setminus u}$

$$= W_{f,S}^I \oplus \bigoplus_{s \in \mathcal{K}(f) \setminus u_f} \bar{W}_{f,S}$$

$$= \begin{cases} W_{f,S}^Q & \text{if } |\mathcal{K}(f) \setminus u_f| \text{ is odd} \\ W_{f,S}^I & \text{otherwise} \end{cases}$$

and  $\Upsilon_{f,S \setminus u_f}^Q \oplus \bar{C}_{f,S \setminus u_f} \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} C_{f,S \setminus u}^Q$

$$= \bar{W}_{f,S} \oplus \bigoplus_{s \in \mathcal{K}(f) \setminus u_f} W_{f,S}^Q$$

$$= \begin{cases} W_{f,S}^I & \text{if } |\mathcal{K}(f) \setminus u_f| \text{ is odd} \\ \bar{W}_{f,S} & \text{otherwise.} \end{cases}$$

If  $|\mathcal{K}(f) \setminus u_f|$  is even, we obtain  $W_{f,S}^Q$  as  $W_{f,S}^Q = \bar{W}_{f,S} \oplus W_{f,S}^I$ . Finally, we have  $W_{f,S} = \langle W_{f,S}^I, W_{f,S}^Q \rangle$ .

4) *Decoding of requested subfiles not cached*: Next, we show how user  $U_k$  obtains all the remaining requested subfiles  $W_{\mathbf{d}(k),\mathcal{S}}$  for all sets  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$ . These subfiles satisfy  $k \notin \mathcal{S}$  and are, thus, not coded in the cache of user  $U_k$ . Hereafter  $f = \mathbf{d}(k)$ . Let us first suppose user  $U_k$  is leader, i.e.  $k = u_f$ . Observe that then, associated to any set  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$ , since  $k = u_f \notin \mathcal{S}$ , there exists a broadcasted subfile  $Y_{f,\mathcal{S}}$ . In addition, user  $k = u_f$  has already obtained all the requested subfiles coded in its cache, in particular subfiles  $W_{f,u_f \cup \mathcal{S} \setminus s}$  for all  $s \in \mathcal{S} \cap \mathcal{K}(f)$ , and can compute  $W_{f,\mathcal{S}} = \langle W_{f,\mathcal{S}}^I, W_{f,\mathcal{S}}^Q \rangle$  as

$$W_{f,\mathcal{S}}^I = Y_{f,\mathcal{S}}^I \oplus \bigoplus_{s \in \mathcal{S} \cap \mathcal{K}(f)} \bar{W}_{f,u_f \cup \mathcal{S} \setminus s}, \quad (25)$$

$$W_{f,\mathcal{S}}^Q = W_{f,\mathcal{S}}^I \oplus Y_{f,\mathcal{S}}^Q \oplus \bigoplus_{s \in \mathcal{S} \cap \mathcal{K}(f)} W_{f,u_f \cup \mathcal{S} \setminus s}^Q. \quad (26)$$

for all sets  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$ . Instead, if user  $U_k$  is no leader, i.e.  $k \in \mathcal{K}(f) \setminus u_f$  he obtains, first, the subfiles  $W_{f,\mathcal{S}}$  for all the set  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$  satisfying  $u_f \in \mathcal{S}$ . Observe that for each of these sets, associated to the set  $\mathcal{S} \cup k \setminus u_f$ , there exists the broadcasted subfile  $Y_{f,\mathcal{S} \cup k \setminus u_f}$  with half-subfiles given by  $Y_{f,\mathcal{S} \cup k \setminus u_f}^I$

$$\begin{aligned} &= W_{f,\mathcal{S} \cup k \setminus u_f}^I \oplus \bigoplus_{s \in \{\mathcal{S} \cup k \setminus u_f\} \cap \mathcal{K}(f)} \bar{W}_{f,u_f \cup \{\mathcal{S} \cup k \setminus u_f\} \setminus s} \\ &= W_{f,\mathcal{S} \cup k \setminus u_f}^I \oplus \bar{W}_{f,\mathcal{S}} \oplus \bigoplus_{s \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} \bar{W}_{f,\mathcal{S} \cup k \setminus s} \end{aligned}$$

and  $Y_{f,\mathcal{S} \cup k \setminus u_f}^Q$

$$= \bar{W}_{f,\mathcal{S} \cup k \setminus u_f} \oplus W_{f,\mathcal{S}}^Q \oplus \bigoplus_{s \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} W_{f,\mathcal{S} \cup k \setminus s}^Q.$$

Given that user  $U_k$  has already obtained all the requested subfiles coded in its own cache, in particular subfiles  $W_{f,\mathcal{S} \cup k \setminus s}$  for all  $s \in \mathcal{S} \cap \mathcal{K}(f)$ , he can compute  $W_{f,\mathcal{S}} = \langle W_{f,\mathcal{S}}^I, W_{f,\mathcal{S}}^Q \rangle$  as

$$W_{f,\mathcal{S}}^Q = Y_{f,\mathcal{S} \cup k \setminus u_f}^Q \oplus \bar{W}_{f,\mathcal{S} \cup k \setminus u_f} \oplus \bigoplus_{s \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} W_{f,\mathcal{S} \cup k \setminus s}^Q \quad (27)$$

$$W_{f,\mathcal{S}}^I = W_{f,\mathcal{S}}^Q \oplus Y_{f,\mathcal{S} \cup k \setminus u_f}^I \oplus W_{f,\mathcal{S} \cup k \setminus u_f}^I \oplus \bigoplus_{s \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} \bar{W}_{f,\mathcal{S} \cup k \setminus s} \quad (28)$$

for all sets  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$  with  $u_f \in \mathcal{S}$ . Now user  $U_k$  can mimic the computations performed at the user leader  $u_f$  in (25) and (26), to obtain the remaining subfiles, i.e.  $W_{f,\mathcal{S}}$  for all sets  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$  satisfying  $u_f \notin \mathcal{S}$ .

### C. Extension to arbitrary requested files

The strategy presented in previous subsection requires all  $N$  files to be requested by at least one user. In this section, we show that the same scheme can be applied for any request if configured to  $K + 1$  users instead of  $K$ .

1) *Prefetching scheme*: We use the subfile partitions and prefetching scheme specified in previous section, but as if there were  $K + 1$  users instead of  $K$ . As a result, each file is partitioned into  $\binom{K+1}{r}$  subfiles and thus, the required cache load at users equals  $MF$ , with

$$M = \frac{r}{K+1} + (N-1) \frac{r}{K+1} \frac{r-1}{K}$$

2) *Delivery scheme*: For the delivery, if all files are requested, since the caching system is configured for  $K + 1$  users, we need to add one file request for the  $K + 1$ th “virtual” user. Given that all files are already requested by some user, the  $K + 1$ th user can request any of the files in  $\mathcal{F}$ . Then, we apply the delivery strategy as described in previous subsections. In that case, every user obtains their desired file, with a rate

$$R = N \left( 1 - \frac{r}{K+1} \right). \quad (29)$$

Instead, if some files are not requested, i.e.  $N_e(\mathbf{d}) \leq N - 1$ , we first transform the coded caching system with  $N$  files into caching system with  $N_e(\mathbf{d}) + 1$  files. To that end, we define a new file  $W_0$  as the binary sum of all files not requested

$$W_0 = \bigoplus_{f \in \mathcal{F} \setminus \mathcal{N}_e(\mathbf{d})} W_f$$

where  $\mathcal{N}_e(\mathbf{d}) \subseteq \mathcal{F}$ . Then, the set of files in the system  $\mathcal{F}$  is given by  $\mathcal{F}_0(\mathbf{d}) = \{0 \cup \mathcal{N}_e(\mathbf{d})\}$ . The server can obtain the subfiles of the new file  $W_0$  from the not requested subfiles as

$$W_{0,\mathcal{S}} = \bigoplus_{f \in \mathcal{F} \setminus \mathcal{N}_e(\mathbf{d})} W_{f,\mathcal{S}}.$$

Similarly, users can consider the coded cached subfiles in (5), as if the set of file was  $\mathcal{F}_0(\mathbf{d})$ , since

$$\begin{aligned} Z_{\mathcal{S}} &= \bigoplus_{f \in \mathcal{F}} W_{f,\mathcal{S}} \\ &= \bigoplus_{f \in \mathcal{F} \setminus \mathcal{N}_e(\mathbf{d})} W_{f,\mathcal{S}} \oplus \bigoplus_{f \in \mathcal{N}_e(\mathbf{d})} W_{f,\mathcal{S}} \\ &= \bigoplus_{f \in 0 \cup \mathcal{N}_e(\mathbf{d})} W_{f,\mathcal{S}}. \end{aligned}$$

and, can obtain the coded subfiles in (6) associated to file  $W_0$  as

$$\begin{aligned} Z_{0,\mathcal{S}^-} &= \bigoplus_{f \in \mathcal{F} \setminus \mathcal{N}_e(\mathbf{d})} Z_{f,\mathcal{S}^-} \\ &= \bigoplus_{s \in \mathcal{K} \setminus \mathcal{S}^-} \bigoplus_{f \in \mathcal{F} \setminus \mathcal{N}_e(\mathbf{d})} W_{f,\mathcal{S}^- \cup s} \\ &= \bigoplus_{s \in \mathcal{K} \setminus \mathcal{S}^-} W_{0,\mathcal{S}}. \end{aligned}$$

Finally, since we have replaced the set of files  $\mathcal{F}$  by  $\mathcal{F}_0(\mathbf{d})$ , we can assume that the  $K + 1$  virtual user requests file  $W_0$ . Then, observe that all files in  $\mathcal{F}_0(\mathbf{d})$  are requested, and thus we can apply the delivery strategy as described in previous subsection, with a rate

$$R = (N_e(\mathbf{d}) + 1) \left( 1 - \frac{r}{K+1} \right). \quad (30)$$

Combining (29) and (30), we can write the rate as presented in Theorem 1.



Coded cached subfiles at user $U_1$	
$Z_{\{1,2,3\}}$	$W_{1,\{1,2,3\}} \oplus W_{2,\{1,2,3\}}$
$Z_{\{1,2,4\}}$	$W_{1,\{1,2,4\}} \oplus W_{2,\{1,2,4\}}$
$Z_{\{1,2,5\}}$	$W_{1,\{1,2,5\}} \oplus W_{2,\{1,2,5\}}$
$Z_{\{1,3,4\}}$	$W_{1,\{1,3,4\}} \oplus W_{2,\{1,3,4\}}$
$Z_{\{1,3,5\}}$	$W_{1,\{1,3,5\}} \oplus W_{2,\{1,3,5\}}$
$Z_{\{1,4,5\}}$	$W_{1,\{1,4,5\}} \oplus W_{2,\{1,4,5\}}$
$Z_{1,\{1,2\}}$	$W_{1,\{1,2,3\}} \oplus W_{1,\{1,2,4\}} \oplus W_{1,\{1,2,5\}}$
$Z_{1,\{1,3\}}$	$W_{1,\{1,3,2\}} \oplus W_{1,\{1,3,4\}} \oplus W_{1,\{1,3,5\}}$
$Z_{1,\{1,4\}}$	$W_{1,\{1,4,2\}} \oplus W_{1,\{1,4,3\}} \oplus W_{1,\{1,4,5\}}$
Coded subfiles computed at user $U_1$	
$Z_{1,\{1,5\}}$	$W_{1,\{1,5,2\}} \oplus W_{1,\{1,5,3\}} \oplus W_{1,\{1,5,4\}}$
$Z_{2,\{1,2\}}$	$W_{2,\{1,2,3\}} \oplus W_{2,\{1,2,4\}} \oplus W_{2,\{1,2,5\}}$
$Z_{2,\{1,3\}}$	$W_{2,\{1,3,2\}} \oplus W_{2,\{1,3,4\}} \oplus W_{2,\{1,3,5\}}$
$Z_{2,\{1,4\}}$	$W_{2,\{1,4,2\}} \oplus W_{2,\{1,4,3\}} \oplus W_{2,\{1,4,5\}}$
$Z_{2,\{1,5\}}$	$W_{2,\{1,5,2\}} \oplus W_{2,\{1,5,3\}} \oplus W_{2,\{1,5,4\}}$

Table I: Coded cached and computable coded subfiles at user  $U_1$  for the caching system of Example 1.

## V. EXAMPLE

In this section, we apply the caching scheme presented in section IV-B to a particular example. Consider a caching system with  $N = 2$  files,  $K = 5$  users, and a caching capacity of  $MF$  bits with  $M = \frac{9}{10}$ , which corresponds to  $r = 3$ . Each file  $W_f$ ,  $f = \{1, 2\}$  is divide into  $\binom{K}{r} = 10$  subfiles, i.e.  $W_{f,\{1,2,3\}}$ ,  $W_{f,\{1,2,4\}}$ ,  $W_{f,\{1,2,5\}}$ ,  $W_{f,\{1,3,4\}}$ ,  $W_{f,\{1,3,5\}}$ ,  $W_{f,\{1,4,5\}}$ ,  $W_{f,\{2,3,4\}}$ ,  $W_{f,\{2,3,5\}}$ ,  $W_{f,\{2,4,5\}}$ , and  $W_{f,\{3,4,5\}}$ .

The cache of user  $U_k$  stores a linear combination of the subfiles  $W_{f,S}$  of both files satisfying  $k \in \mathcal{S}$ . For the running example, the cache of user  $U_1$  is detailed in Table I. Observe that in this case, we have chosen to omit coded cached subfiles  $Z_{f,S^-}$  satisfying  $f = 2$  or  $5 \in \mathcal{S}^-$ . These subfiles (see the Table I) can be computed from the cache content at user  $U_1$  as

$$\begin{aligned}
Z_{1,\{1,5\}} &= Z_{1,\{1,2\}} \oplus Z_{1,\{1,3\}} \oplus Z_{1,\{1,4\}}, \\
Z_{2,\{1,2\}} &= Z_{1,\{1,2\}} \oplus Z_{\{1,2,3\}} \oplus Z_{\{1,2,4\}} \oplus Z_{\{1,2,5\}}, \\
Z_{2,\{1,3\}} &= Z_{1,\{1,3\}} \oplus Z_{\{1,3,2\}} \oplus Z_{\{1,3,4\}} \oplus Z_{\{1,3,5\}}, \\
Z_{2,\{1,4\}} &= Z_{1,\{1,4\}} \oplus Z_{\{1,4,2\}} \oplus Z_{\{1,4,3\}} \oplus Z_{\{1,4,5\}}, \\
Z_{2,\{1,5\}} &= Z_{1,\{1,5\}} \oplus Z_{\{1,5,2\}} \oplus Z_{\{1,5,3\}} \oplus Z_{\{1,5,4\}}.
\end{aligned}$$

Given that every file is divided into 10 subfiles, and each user caches 9 coded subfiles, the cache load condition  $M = \frac{9}{10}$  is satisfied.

Given the above prefetching scheme, we illustrate our proposed delivery strategy for a representative demand, where users  $U_1, U_3, U_5$  request file  $W_1$ , and users  $U_2$  and  $U_4$  request file  $W_2$ . This corresponds to the demand vector  $\mathbf{d} = [1, 2, 1, 2, 1]$ , and  $\mathcal{K}(1) = \{1, 3, 5\}$  and  $\mathcal{K}(2) = \{2, 4\}$ . Observe that all files are requested. First, the server selects one user leader for each file, i.e.  $u_1 = 1$  and  $u_2 = 2$ . Then, the server broadcasts, according to (9) and (10), the coded subfiles  $Y_{f,S} = \langle Y_{f,S}^I, Y_{f,S}^Q \rangle$  specified in Table II. Observe that the total number of broadcasted subfiles is  $N \binom{K-1}{r} = 8$

Next, we detail the decoding process at user  $U_1$ . He first obtains the subfiles  $W_{1,S}$  coded its cache, i.e.  $1 \in \mathcal{S}$ .

$Y_{1,\{2,3,4\}}$	$Y_{1,\{2,3,4\}}^I$	$W_{1,\{2,3,4\}}^I \oplus W_{1,\{1,2,4\}}$
	$Y_{1,\{2,3,4\}}^Q$	$\bar{W}_{1,\{2,3,4\}} \oplus \bar{W}_{1,\{1,2,4\}}^Q$
$Y_{1,\{2,3,5\}}$	$Y_{1,\{2,3,5\}}^I$	$W_{1,\{2,3,5\}}^I \oplus W_{1,\{2,1,5\}} \oplus W_{1,\{2,3,1\}}$
	$Y_{1,\{2,3,5\}}^Q$	$\bar{W}_{1,\{2,3,5\}} \oplus \bar{W}_{1,\{2,1,5\}}^Q \oplus \bar{W}_{1,\{2,3,1\}}^Q$
$Y_{1,\{2,4,5\}}$	$Y_{1,\{2,4,5\}}^I$	$W_{1,\{2,4,5\}}^I \oplus W_{1,\{2,4,1\}}$
	$Y_{1,\{2,4,5\}}^Q$	$\bar{W}_{1,\{2,4,5\}} \oplus \bar{W}_{1,\{2,4,1\}}^Q$
$Y_{1,\{3,4,5\}}$	$Y_{1,\{3,4,5\}}^I$	$W_{1,\{3,4,5\}}^I \oplus W_{1,\{1,4,5\}} \oplus W_{1,\{3,4,1\}}$
	$Y_{1,\{3,4,5\}}^Q$	$\bar{W}_{1,\{3,4,5\}} \oplus \bar{W}_{1,\{1,4,5\}}^Q \oplus \bar{W}_{1,\{3,4,1\}}^Q$
$Y_{2,\{1,3,4\}}$	$Y_{2,\{1,3,4\}}^I$	$W_{2,\{1,3,4\}}^I \oplus W_{2,\{1,3,2\}}$
	$Y_{2,\{1,3,4\}}^Q$	$\bar{W}_{2,\{1,3,4\}} \oplus \bar{W}_{2,\{1,3,2\}}^Q$
$Y_{2,\{1,3,5\}}$	$Y_{2,\{1,3,5\}}^I$	$W_{2,\{1,3,5\}}^I$
	$Y_{2,\{1,3,5\}}^Q$	$\bar{W}_{2,\{1,3,5\}}^Q$
$Y_{2,\{1,4,5\}}$	$Y_{2,\{1,4,5\}}^I$	$W_{2,\{1,4,5\}}^I \oplus W_{2,\{1,2,5\}}$
	$Y_{2,\{1,4,5\}}^Q$	$\bar{W}_{2,\{1,4,5\}} \oplus \bar{W}_{2,\{1,2,5\}}^Q$
$Y_{2,\{3,4,5\}}$	$Y_{2,\{3,4,5\}}^I$	$W_{2,\{3,4,5\}}^I \oplus W_{2,\{3,2,5\}}$
	$Y_{2,\{3,4,5\}}^Q$	$\bar{W}_{2,\{3,4,5\}} \oplus \bar{W}_{2,\{3,2,5\}}^Q$

Table II: Broadcasted coded half-subfiles for the demand  $\mathbf{d} = [1, 2, 1, 2, 1]$

$p =  \mathcal{S} \cap \mathcal{K}(2) $	$u_2 = 2$	$W_{2,S}$
0	$u_2 \notin \mathcal{S}$	$W_{2,\{1,3,5\}}$
1	$u_2 \in \mathcal{S}$	$W_{2,\{1,2,3\}}$
		$W_{2,\{1,2,5\}}$
	$u_2 \notin \mathcal{S}$	$W_{2,\{1,3,4\}}$
2	$u_2 \in \mathcal{S}$	$W_{2,\{1,2,4\}}$

Table III: Classification of interfering subfiles  $W_{2,S}$  at user  $U_1$ , i.e.  $1 \in \mathcal{S}$

According to (11), user  $U_1$  can obtain these subfiles by, first, obtaining the interfering subfile  $W_{2,S}$ , and then computing

$$W_{1,S} = Z_S \oplus W_{2,S}$$

for all sets  $\mathcal{S}$  satisfying  $1 \in \mathcal{S}$ . To that end, user  $U_1$  divides the decoding of these interfering subfiles into 3 consecutive phases  $p = \{0, 1, 2\}$ . In phase  $p$ , user  $U_1$  obtains all subfiles  $W_{2,S}$  satisfying  $1 \in \mathcal{S}$  and  $p = |\mathcal{S} \cap \mathcal{K}(2)|$ . Moreover, within each phase he obtains, first, the subfiles that are coded cached at the user leader,  $u_2 = 2$ , i.e.  $2 \in \mathcal{S}$ . For the running example, the interfering subfiles at user  $U_1$  are classified in Table III according to the above decoding phases. The only phase 0 interfering subfile is  $W_{2,\{1,3,5\}}$ , which can be obtained according to (12) as

$$\begin{aligned}
W_{2,\{1,3,5\}}^I &= Y_{2,\{1,3,5\}}^I, \\
W_{2,\{1,3,5\}}^Q &= Y_{2,\{1,3,5\}}^Q \oplus Y_{2,\{1,3,5\}}^I.
\end{aligned}$$

Phase 1 interfering subfiles satisfying  $u_2 \in \mathcal{S}$  are  $W_{2,\{1,2,3\}}$  and  $W_{2,\{1,2,5\}}$ . Observe that for both sets  $\mathcal{S} \in \{\{1, 2, 3\}, \{1, 2, 5\}\}$ , from (17) and (18), we obtain

$$\begin{aligned}
\Upsilon_{2,S}^I &= \Upsilon_{2,S \setminus 2}^I \oplus C_{2,S \setminus 2}^I, \\
\Upsilon_{2,S}^Q &= \Upsilon_{2,S \setminus 2}^Q \oplus \bar{C}_{2,S \setminus 2}
\end{aligned}$$

and from (15) and (16), respectively

$$\begin{aligned}
C_{2,S \setminus 2} &= Z_{2,S \setminus 2} \oplus W_{2,\{1,3,5\}}, \\
\Upsilon_{2,S \setminus 2} &= Y_{2,S \setminus 2}.
\end{aligned}$$

Then, given that  $|\mathcal{K}(2) \setminus u_f| = 1$  is odd, user  $U_1$  obtains

$$W_{2,S} = \langle W_{2,S}^I, W_{2,S}^Q \rangle \text{ as}$$

$$\begin{aligned} W_{2,S}^I &= \Upsilon_{2,S}^Q \\ &= Y_{2,\{S \setminus 2\} \cup 4}^Q \oplus [\bar{Z}_{2,S \setminus 2} \oplus \bar{W}_{2,\{1,3,5\}}], \end{aligned}$$

$$\begin{aligned} W_{2,S}^Q &= \Upsilon_{2,S}^I \\ &= Y_{2,\{S \setminus 2\} \cup 4}^I \oplus [Z_{2,S \setminus 2}^I \oplus W_{2,\{1,3,5\}}^I]. \end{aligned}$$

Phase 1 interfering subfiles satisfying  $u_2 \notin S$  are  $W_{2,\{1,3,4\}}$  and  $W_{2,\{1,4,5\}}$ . User  $U_1$  obtains this interfering subfile according to (13) and (14) as

$$\begin{aligned} W_{2,\{1,3,4\}}^I &= Y_{2,\{1,3,4\}}^I \oplus \bar{W}_{2,\{1,2,3\}}, \\ W_{2,\{1,3,4\}}^Q &= Y_{2,\{1,3,4\}}^Q \oplus W_{2,\{1,2,3\}}^Q \oplus W_{2,\{1,3,4\}}^I, \end{aligned}$$

and

$$\begin{aligned} W_{2,\{1,4,5\}}^I &= Y_{2,\{1,4,5\}}^I \oplus \bar{W}_{2,\{1,2,5\}}, \\ W_{2,\{1,4,5\}}^Q &= Y_{2,\{1,4,5\}}^Q \oplus W_{2,\{1,2,5\}}^Q \oplus W_{2,\{1,4,5\}}^I. \end{aligned}$$

Finally, user  $U_1$  can obtain the remaining phase 2 subfile  $W_{2,\{1,2,4\}}$ , according to (17) and (18). However, given that this is the only interfering subfile left, he can obtain it directly from, either  $Z_{2,\{1,4\}}$  or  $Z_{2,\{1,2\}}$ , as

$$\begin{aligned} W_{2,\{1,2,4\}} &= Z_{2,\{1,4\}} \oplus W_{2,\{1,4,3\}} \oplus W_{2,\{1,4,5\}}, \\ W_{2,\{1,2,4\}} &= Z_{2,\{1,2\}} \oplus W_{2,\{1,2,3\}} \oplus W_{2,\{1,2,5\}}. \end{aligned}$$

Once all interfering subfiles  $W_{2,S}$  for all  $S$  satisfying  $1 \in S$  are decoded. User  $U_1$  decodes the requested subfiles only available at other users' cache. These are the subfiles  $W_{1,S}$  for  $S \in \{\{2, 3, 4\}, \{2, 3, 5\}, \{2, 4, 5\}\}$ . Given that user  $U_1$  is leader, he obtains these subfiles according to (25) and (26), as

$$W_{1,\{2,3,4\}}^I = Y_{1,\{2,3,4\}}^I \oplus \bar{W}_{1,\{1,2,4\}}, \quad (31)$$

$$W_{1,\{2,3,4\}}^Q = W_{1,\{2,3,4\}}^I \oplus Y_{1,\{2,3,4\}}^Q \oplus W_{1,\{1,2,4\}}^Q, \quad (32)$$

$$W_{1,\{2,3,5\}}^I = Y_{1,\{2,3,5\}}^I \oplus \bar{W}_{1,\{1,2,5\}} \oplus \bar{W}_{1,\{1,2,3\}}, \quad (33)$$

$$W_{1,\{2,3,5\}}^Q = W_{1,\{2,3,5\}}^I \oplus Y_{1,\{2,3,5\}}^Q \oplus W_{1,\{1,2,5\}}^Q \oplus W_{1,\{1,2,3\}}^Q, \quad (34)$$

$$W_{1,\{2,4,5\}}^I = Y_{1,\{2,4,5\}}^I \oplus \bar{W}_{1,\{1,2,4\}}, \quad (35)$$

$$W_{1,\{2,4,5\}}^Q = W_{1,\{2,4,5\}}^I \oplus Y_{1,\{2,4,5\}}^Q \oplus W_{1,\{1,2,4\}}^Q. \quad (36)$$

For completeness, as the decoding of the requested subfiles at other users' cache depends on whether the decoding user is leader or not, let us also consider the decoding process at the non user leader  $U_3$ . User  $U_3$  first obtains the subfiles  $W_{1,\{1,2,4\}}$ ,  $W_{1,\{1,2,5\}}$ ,  $W_{1,\{1,4,5\}}$  coded cached at the user leader  $U_1$ , according to (27) and (28) as

$$W_{1,\{1,2,4\}}^Q = Y_{1,\{2,3,4\}}^Q \oplus \bar{W}_{1,\{2,3,4\}},$$

$$W_{1,\{1,2,4\}}^I = W_{1,\{1,2,4\}}^Q \oplus Y_{1,\{2,3,4\}}^I \oplus W_{1,\{2,3,4\}}^I,$$

$$W_{1,\{1,2,5\}}^Q = Y_{1,\{2,3,5\}}^Q \oplus \bar{W}_{1,\{2,3,5\}} \oplus W_{1,\{1,2,3\}}^Q,$$

$$W_{1,\{1,2,5\}}^I = W_{1,\{1,2,5\}}^Q \oplus Y_{1,\{2,3,5\}}^I \oplus W_{1,\{2,3,5\}}^I \oplus \bar{W}_{1,\{1,2,3\}},$$

$$W_{1,\{1,4,5\}}^Q = Y_{1,\{3,4,5\}}^Q \oplus \bar{W}_{1,\{3,4,5\}} \oplus W_{1,\{1,3,4\}}^Q,$$

$$W_{1,\{1,4,5\}}^I = W_{1,\{1,4,5\}}^Q \oplus Y_{1,\{3,4,5\}}^I \oplus W_{1,\{3,4,5\}}^I \oplus \bar{W}_{1,\{1,3,4\}}.$$

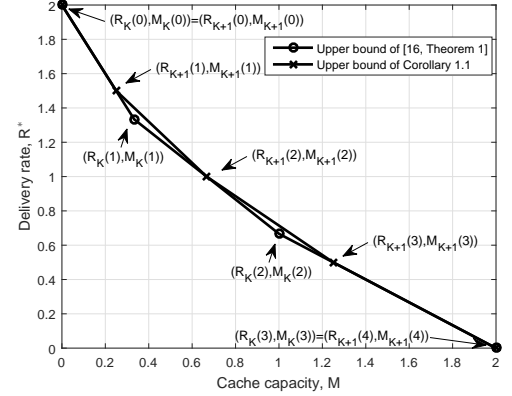


Figure 2: Rate-memory functions of the proposed scheme and the Tian-Chen scheme for  $N = 2$  and  $K = 3$ .

Finally, user  $U_3$  obtains subfile  $W_{1,\{2,4,5\}}$ , as in (35) and (36).

## VI. CONCLUSIONS

In this work, we proposed a novel caching scheme that approaches the rate-memory region achieved by the Tian-Chen scheme as the number of users in the system increases, which only requires a finite field of order  $2^2$ . Moreover, instead of relying on the existence of some valid code, we provided an explicit combinatorial construction of the caching scheme.

## APPENDIX A

### PROOF OF COROLLARY 1.2

Let us define

$$R_K^*(r) \triangleq N \left(1 - \frac{r}{K}\right),$$

$$M_K(r) \triangleq \frac{r}{K} + (N-1) \frac{r(r-1)}{K(K-1)}$$

for all  $r \in \{0, \dots, K\}$ . In a caching system with  $K$  users and  $N$  files, the worst demand rate-memory pairs, as a function of the parameter  $r$ , for the Tian-Chen scheme are given by  $R_{TC}(r) = R_K(r)$  and  $M_{TC}(r) = M_K(r)$ , and for the strategy presented here by  $R_{TC+}^*(r) = R_{K+1}(r)$  and  $M_{TC+}(r) = M_{K+1}(r)$ .

To support this proof, in Fig. 2 we compare the rate-memory functions obtained by both schemes in a particular caching system with  $N = 2$  files and  $K = 3$  user. We first obtain the memory points and intervals where the Tian-Chen and the proposed strategy coincide. By memory sharing between the rate-memory pairs in  $(R_K^*(r), M_K(r))$  at  $r$  and  $r+1$ , we obtain the rate-memory pairs

$$\begin{aligned} R_K^*(r, y) &= (1-y)R_K^*(r) + yR_K^*(r+1) \\ &= N \left(1 - \frac{r}{K}\right) - y \frac{N}{K} \end{aligned} \quad (37)$$

$$\begin{aligned} M_K(r, y) &= (1-y)M_K(r) + yM_K(r+1) \\ &= \frac{r}{K} + (N-1) \frac{r(r-1)}{K(K-1)} \\ &\quad + y \left( \frac{1}{K} + 2 \frac{r(N-1)}{K(K-1)} \right) \end{aligned} \quad (38)$$

for  $y \in [0, 1]$  and  $r = \{0, \dots, K-1\}$ . Next, particularizing (37) and (38) at  $\hat{y}(r) = \frac{K-r}{K+1}$ , we obtain

$$\begin{aligned} R_K^*(r, \hat{y}(r)) &= R_{K+1}^*(r+1), \\ M_K(r, \hat{y}(r)) &= M_{K+1}(r+1) \end{aligned}$$

for  $r \in \{0, \dots, K\}$ . This is, the delivery rate for the Tian-Chen and the proposed scheme coincide  $R_{\text{TC}^+}^*(M) = R_{\text{TC}}^*(M)$  at  $M = M_{K+1}(r+1) = \frac{r+1}{K+1} (1 + (N-1) \frac{r}{K})$  for  $r = \{0, \dots, K\}$ . Furthermore, at  $r = 0$ , we have  $R_{K+1}^*(0) = R_K^*(0) = N$ , and  $M_{K+1}(0) = M_K(0) = 0$  and thus, equality holds also for the memory interval

$$0 \leq M \leq M_{K+1}(1) = \frac{1}{K+1}.$$

Similarly, we have  $R_{K+1}^*(K+1) = R_K^*(K) = 0$ , and  $M_{K+1}(K+1) = M_K(K) = N$ , and thus, the equality holds also for the memory interval

$$N \geq M \geq \frac{K}{K+1} + (N-1) \frac{K-1}{K+1}.$$

Next, we obtain the bound on the maximum rate ratio between both strategies. It can be shown that the equality  $M_K(r) = M_{K+1}(r, y)$  is achieved at  $y = \check{y}(r)$  with

$$\check{y}(r) = \frac{r}{K-1} \left( 1 - \frac{2(N-1)+1}{K+2(N-1)r} \right)$$

for  $r = \{1, \dots, K-1\}$ . Given that both rate-memory functions coincide at  $M = M_{K+1}(r+1)$ , that  $M_{K+1}(r) \leq M_K(r) \leq M_{K+1}(r+1)$  for  $r = \{1, \dots, K-1\}$ , and the convexity of the two piece-wise linear rate-memory functions, we have, see Fig. 2, that the maximum rate ratio must be found at one of the memory points  $M_K(r) = M_{K+1}(r, \check{y}(r))$ ,  $r = \{1, \dots, K-1\}$ , this is

$$\begin{aligned} \max_{0 \leq M \leq N} \frac{R^*(M)}{R_{\text{TC}}^*(M)} &= \max_{r \in \{0, \dots, K-1\}} \frac{R_{K+1}(r, \check{y})}{R_K(r)} \\ &= 1 + \frac{r - \check{y}(r)K}{(K-r)(K+1)} \\ &\leq 1 + \frac{1}{(K-1)(K+1)}. \end{aligned} \quad (39)$$

To obtain inequality (39), observe that  $\check{y}(r)$  decreases as a function of  $N$ , and thus, as  $N$  approaches infinity, we have for  $r = \{1, \dots, K-1\}$   $\lim_{N \rightarrow \infty} \check{y}(r) = \frac{r}{K-1} (1 - \frac{1}{r})$ .

## 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, May 2014.
- [2] M. R. Korupolu, C. Plaxton, and R. Rajaraman, "Placement algorithms for hierarchical cooperative caching," *Journal of Algorithms*, vol. 38, no. 1, pp. 260–302, 2001.
- [3] A. Dan, D. Sitaram, and P. Shahabuddin, "Dynamic batching policies for an on-demand video server," *Multimedia Systems*, vol. 4, no. 3, pp. 112–121, 1996.
- [4] 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, Aug 2015.
- [5] U. Niesen and M. A. Maddah-Ali, "Coded caching with nonuniform demands," *IEEE Transactions on Information Theory*, vol. 63, no. 2, pp. 1146–1158, Feb 2017.
- [6] J. Zhang, X. Lin, and X. Wang, "Coded caching under arbitrary popularity distributions," in *ITA workshop*, Feb 2015, pp. 98–107.

- [7] M. Ji, A. M. Tulino, J. Llorca, and G. Caire, "Order-optimal rate of caching and coded multicasting with random demands," *IEEE Transactions on Information Theory*, vol. 63, no. 6, pp. 3923–3949, June 2017.
- [8] R. Pedarsani, M. A. Maddah-Ali, and U. Niesen, "Online coded caching," *IEEE/ACM Transactions on Networking*, vol. 24, no. 2, pp. 836–845, April 2016.
- [9] 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, June 2016.
- [10] M. M. Amiri, Q. Yang, and D. Gündüz, "Coded caching for a large number of users," in *2016 IEEE Information Theory Workshop (ITW)*, Sept 2016, pp. 171–175.
- [11] S. Sahraei and M. Gastpar, "K users caching two files: An improved achievable rate," in *2016 Annual Conference on Information Science and Systems (CISS)*, March 2016, pp. 620–624.
- [12] H. Ghasemi and A. Ramamoorthy, "Improved lower bounds for coded caching," in *IEEE ISIT*, June 2015, pp. 1696–1700.
- [13] M. M. Amiri and D. Gündüz, "Fundamental limits of coded caching: Improved delivery rate-cache capacity tradeoff," *IEEE Transactions on Communications*, vol. 65, no. 2, pp. 806–815, Feb 2017.
- [14] Z. Chen, P. Fan, and K. B. Letaief, "Fundamental limits of caching: improved bounds for users with small buffers," *IET Communications*, vol. 10, no. 17, pp. 2315–2318, 2016.
- [15] Q. Yu, M. A. Maddah-Ali, and A. S. Avestimehr, "The exact rate-memory tradeoff for caching with uncoded prefetching," *IEEE Transactions on Information Theory*, vol. 64, no. 2, pp. 1281–1296, Feb 2018.
- [16] C. Tian and J. Chen, "Caching and delivery via interference elimination," *IEEE Transactions on Information Theory*, vol. 64, no. 3, pp. 1548–1560, March 2018.
- [17] C. Tian and K. Zhang, "From uncoded prefetching to coded prefetching in coded caching," *CoRR*, vol. abs/1704.07901, 2017. [Online]. Available: <http://arxiv.org/abs/1704.07901>
- [18] K. Wan, D. Tuninetti, and P. Piantanida, "On caching with more users than files," in *IEEE ISIT*, July 2016, pp. 135–139.
- [19] A. Sengupta, R. Tandon, and T. C. Clancy, "Improved approximation of storage-rate tradeoff for caching via new outer bounds," in *IEEE ISIT*, June 2015, pp. 1691–1695.
- [20] C. Tian, "A note on the fundamental limits of coded caching," *CoRR*, vol. abs/1503.00010, 2015.
- [21] 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 ISIT*, June 2017, pp. 386–390.
- [22] C. Y. Wang, S. S. Bidokhti, and M. Wigger, "Improved converses and gap-results for coded caching," pp. 2428–2432, June 2017.
- [23] K. Shanmugam, A. M. Tulino, and A. G. Dimakis, "Coded caching with linear subpacketization is possible using ruzsa-szemerédi graphs," in *2017 IEEE International Symposium on Information Theory (ISIT)*, June 2017, pp. 1237–1241.
- [24] L. Tang and A. Ramamoorthy, "Low subpacketization schemes for coded caching," in *2017 IEEE International Symposium on Information Theory (ISIT)*, June 2017, pp. 2790–2794.
- [25] K. Shanmugam, M. Ji, A. M. Tulino, J. Llorca, and A. G. Dimakis, "Finite-length analysis of caching-aided coded multicasting," *IEEE Transactions on Information Theory*, vol. 62, no. 10, pp. 5524–5537, Oct 2016.
- [26] U. Niesen and M. A. Maddah-Ali, "Coded caching for delay-sensitive content," in *2015 IEEE International Conference on Communications (ICC)*, June 2015, pp. 5559–5564.
- [27] J. Gómez-Vilardebó, "Fundamental Limits of Caching: Improved Bounds with Coded Prefetching," *ArXiv e-prints*, Dec. 2016.



**Jesús Gómez-Vilardebó** received his M.Sc. and Ph.D. degrees in Telecommunication Engineering from the Universitat Politècnica de Catalunya (UPC) in October 2003 and July 2009, respectively. In September 2005, he was granted by the Centre Tecnològic de Telecomunicacions de Catalunya (CTTC/CERCA) to obtain the Ph.D. in Signal Theory and Communications at the UPC. He is now with the CTTC holding a Research Associate position. His research interests are in the areas of information and communication theory, and their applications in multi-user communications, coded caching and information privacy.