

# A novel coded caching scheme with coded prefetching

Jesus Gomez-Vilardebo

Centre Tecnològic de Telecomunicacions de Catalunya (CTTC/CERCA)

Castelldefels, Spain

jesus.gomez@cttc.es

**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  $m \geq K \log_2(N)$  in some situations, 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

The 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 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, there 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.

This work was partially supported by the Catalan Government under grant SGR2017-1479 and the Spanish Government under grant TEC2013-44591-P (INTENSIV).

The fundamental caching scheme developed in [1] was later extended to more realistic situations. The decentralized setting was considered in [2], non-uniform demands in [3], online coded caching in [4], and hierarchical cache network in [5], among others. In addition, new schemes pushing further the fundamental limits of caching systems have appeared in [6]–[12]. 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 [13]. Other lower bounds have appeared in [7], [14], [15].

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 [11]. 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 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 [12], 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. Finally, conclusions are drawn in Section V.

## 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  $N \leq K$ . As 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. 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. 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 demand  $\mathbf{d}$ , we say that communication rate  $R(\mathbf{d}, \mathcal{M})$  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  $W_{\mathbf{d}(k)}$ . 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 worst demand rate-memory pair  $(R^*, M)$  and the worst demand rate-memory function  $R^*(M)$ .

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 [16, Corollari 1.1], the rate-memory pairs obtained in [10, Corollary 1], and the rate-memory pairs obtained in [11, Theorem 1]

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

with  $r \in \{0, 1, \dots, K\}$ . The rate-memory pairs in [10, Corollary 1] include as special case the rate-memory pair obtained in [6] for  $M = \frac{N}{K}$ . In addition, both [16, Corollari 1.1], and [11, Theorem 1] include the rate-memory pair  $(R^*, M) = (N - \frac{N}{K}, \frac{1}{K})$  obtained in [9]. Recently, a method to obtain new rate-memory pairs has been developed in [12]. However, no explicit characterization of the achievable rate-memory pairs is provided. While the schemes in [10] and [16] only require binary operations, the scheme developed in [11], [12] makes use of codes in a finite field of order  $N^K$  in some situations. Finally, there have also been efforts to obtain lower bounds on the delivery rate. The cut-set bound was studied in [1]. A tighter lower bound was obtained in [13]. Other lower bounds have appeared in [7], [14], [15].

### 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

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

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

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

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.

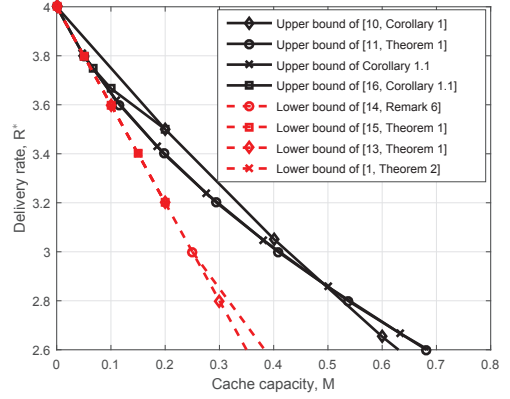


Figure 1: Rate-memory functions of the proposed scheme compared with existing schemes and lower bounds in the literature for  $N = 4$  and  $K = 20$ .

*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  $(R^*, M) = (R_{\text{TC}^+}^*, M_{\text{TC}^+})$ , with

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

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

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

*Remark 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 (1) for  $K+1$  users. The following corollary compares both rate-memory functions. The proof is omitted due to space limitations it can be found in [17].

*Corollary 1.2.* Let  $R_{\text{TC}^+}^*(M)$  be the worst demand rate-memory function achieved by memory sharing between the rate-memory pairs in Corollary 1.1, and let  $R_{\text{TC}}^*(M)$  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  $\frac{R_{\text{TC}^+}^*(M)}{R_{\text{TC}}^*(M)} \leq 1 + \frac{1}{(K-1)(K+1)}$ .

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 [10, Corollary 1], [16, Corollary 1.1] and [11, Theorem 1]. In addition, we depict for comparison, the cut set lower bound [1, Theorem 2], the information-theoretical lower bound obtained in [13, Theorem 1], and the lower bounds recently appeared in [15, Theorem 1] and [14, Remark 6]. Observe that, for the special situation considered here, the new proposed scheme obtains a rate-memory function overlapping almost exactly the rate-memory function in [11], [12]. Moreover, the strategy presented here only requires a finite field of order  $2^2$ . Instead, to obtain the rate-memory pair associated to  $r = 11$ , the Tian-Chen scheme requires a finite field of at least order  $10^5$ .

#### IV. PROPOSED CACHING SCHEME

We prove Theorem 1 in the following subsections. First, we describe a caching scheme that achieves the Tian-Chen rate-memory pairs only if all files are requested by some user. Then, we show that the same scheme can be applied to any demand if configured for  $K + 1$  caches instead of  $K$ .

##### A. 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,  $W_{f,S}$ , one for each subset  $\mathcal{S}$  of  $r$  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 user  $U_k$ . User  $U_k$  caches the binary sum of the  $N$  subfiles associated to  $\mathcal{S}$

$$Z_{\mathcal{S}} = \bigoplus_{f \in \mathcal{F}} W_{f,S} \quad (2)$$

for all  $\mathcal{S} \in \mathcal{T}_k(r)$ . Consequently, each user caches  $\binom{K-1}{r-1}$  coded subfiles. If  $r \geq 2$ , in addition, we need user  $U_k$  to be able to obtain from its cache the coded subfiles

$$Z_{f,S^-} = \bigoplus_{s \in \mathcal{K} \setminus \mathcal{S}^-} W_{f,S^- \cup s} \quad (3)$$

for every file  $f \in \mathcal{F}$ , and every subset  $\mathcal{S}^- \in \mathcal{T}_k(r-1)$ . It can be shown, see [17], that user  $U_k$  only needs to cache the coded subfiles  $Z_{f,S^-}$  in (3) for all files  $f \in \mathcal{F} \setminus g$  and all sets  $\mathcal{S}^- \in \mathcal{T}_k(r-1)$  satisfying  $l \notin \mathcal{S}^-$ , for an arbitrary file  $g$ , and user  $l$ . The total number of coded subfiles  $Z_{f,S^-}$  cached at user  $U_k$  is then  $(N-1)\binom{K-2}{r-2}$ .

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, selects one user leader  $u_f \in \mathcal{K}(f)$  for each file  $f \in \mathcal{F}$ , and then broadcasts the coded subfile  $Y_{f,S} = \langle Y_{f,S}^I, Y_{f,S}^Q \rangle$

$$Y_{f,S}^I = W_{f,S}^I \oplus \bigoplus_{s \in \mathcal{S} \cap \mathcal{K}(f)} \bar{W}_{f,u_f \cup \{s\}}, \quad (4)$$

$$Y_{f,S}^Q = \bar{W}_{f,S} \oplus \bigoplus_{s \in \mathcal{S} \cap \mathcal{K}(f)} W_{f,u_f \cup \{s\}}^Q \quad (5)$$

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  $\mathcal{S} \in \mathcal{T}(r)$  with  $u_f \notin \mathcal{S}$  for each file, we need  $N\binom{K-1}{r}$  subfile length broadcast transmissions, thus, requiring the rate  $R = N\binom{K-1}{r}/\binom{K}{r} = N\left(1 - \frac{r}{K}\right)$ .

Next we describe the decoding process. Each user  $U_k$  divides the decoding of all the requested subfiles  $W_{\mathbf{d}(k),\mathcal{S}}$  for all  $\mathcal{S} \in \mathcal{T}(r)$  into two steps. First, user  $U_k$  obtains all subfiles  $W_{\mathbf{d}(k),\mathcal{S}}$  for all sets  $\mathcal{S} \in \mathcal{T}_k(r)$ . These 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  $W_{\mathbf{d}(k),\mathcal{S}}$  for all sets  $\mathcal{S} \in \bar{\mathcal{T}}_k(r)$ , with  $\bar{\mathcal{T}}_k(r) = \{\mathcal{S} \in \mathcal{T}(r) : k \notin \mathcal{S}\}$ .

3) *Decoding of requested subfiles coded in the cache of user  $U_k$ :* Recall that for any set  $\mathcal{S} \in \mathcal{T}_k(r)$ , user  $U_k$  caches  $Z_{\mathcal{S}} = W_{\mathbf{d}(k),\mathcal{S}} \oplus \bigoplus_{f \in \mathcal{F} \setminus \mathbf{d}(k)} W_{f,S}$  and thus, can obtain  $W_{\mathbf{d}(k),\mathcal{S}}$  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

$$W_{\mathbf{d}(k),\mathcal{S}} = Z_{\mathcal{S}} \oplus \bigoplus_{f \in \mathcal{F} \setminus \mathbf{d}(k)} W_{f,S}. \quad (8)$$

The decoding of these interfering subfiles is divided into  $r$  consecutive phases  $p = 0, \dots, r-1$ . At phase  $p$ , for 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  $W_{f,S}$  as  $\langle W_{f,S}^I, W_{f,S}^Q \rangle = \langle Y_{f,S}^I, Y_{f,S}^Q \oplus Y_{f,S}^I \rangle$ . 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,  $W_{f,S} = \langle W_{f,S}^I, W_{f,S}^Q \rangle$

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

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

Observe that 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  $v \in \bar{\mathcal{K}}(f) \setminus \{\mathcal{S} \setminus u\}$ . This allows user  $U_k$  to compute

$$C_{f,\mathcal{S} \setminus u} = Z_{f,\mathcal{S} \setminus u} \oplus \bigoplus_{v \in \bar{\mathcal{K}}(f) \setminus \{\mathcal{S} \setminus u\}} W_{f,\{\mathcal{S} \setminus u\} \cup v} \quad (11)$$

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

$$\Upsilon_{f,\mathcal{S} \setminus u_f} = \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}. \quad (12)$$

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

$$\Upsilon_{f,S \setminus u_f}^I = C_{f,S \setminus u_f}^I \oplus W_{f,S}^I \oplus \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} \bar{W}_{f,S} \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} \bar{C}_{f,S \setminus u} \oplus \bar{W}_{f,S} \quad (6)$$

$$\Upsilon_{f,S \setminus u_f}^Q = \bar{C}_{f,S \setminus u_f} \oplus \bar{W}_{f,S} \oplus \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} W_{f,S}^Q \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} C_{f,S \setminus u}^Q \oplus W_{f,S}^Q. \quad (7)$$

$\Upsilon_{f,S \setminus u_f}^I \oplus \Upsilon_{f,S \setminus u_f}^Q$ , and  $\bar{C}_{f,S \setminus u} \triangleq C_{f,S \setminus u}^I \oplus C_{f,S \setminus u}^Q$ . Then, user  $U_k$  computes

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

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

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

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

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

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

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

from which we can write

$$\bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} W_{f,\{\mathcal{S} \setminus u\} \cup s} = C_{f,S \setminus u} \oplus W_{f,S}. \quad (13)$$

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

$$\begin{aligned} &= W_{f,\{\mathcal{S} \setminus u_f\} \cup s}^I \oplus \bigoplus_{u \in \{\{\mathcal{S} \setminus u_f\} \cup s\} \cap \mathcal{K}(f)} \bar{W}_{f,u_f \cup \{\{\mathcal{S} \setminus u_f\} \cup s\} \setminus u} \\ &= W_{f,\{\mathcal{S} \setminus u_f\} \cup s}^I \oplus \bar{W}_{f,S} \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} \bar{W}_{f,\{\mathcal{S} \setminus u\} \cup s} \quad (14) \end{aligned}$$

Similarly, we can write  $Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}^Q$

$$= \bar{W}_{f,\{\mathcal{S} \setminus u_f\} \cup s} \oplus W_{f,S}^Q \oplus \bigoplus_{u \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} W_{f,\{\mathcal{S} \setminus u\} \cup s}^Q. \quad (15)$$

Now, it follows from (13) and (14) that  $\Upsilon_{f,S \setminus u_f}^I = \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}^I$  can be written as in (6), and, from (13) and (15) that  $\Upsilon_{f,S \setminus u_f}^Q = \bigoplus_{s \in \mathcal{K}(f) \setminus \mathcal{S}} Y_{f,\{\mathcal{S} \setminus u_f\} \cup s}^Q$  can be written as in (7). Next observe that

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

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

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

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

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

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

$$W_{f,S}^I = Y_{f,S}^I \oplus \bigoplus_{s \in S \cap \mathcal{K}(f)} \bar{W}_{f,u_f \cup S \setminus s}, \quad (16)$$

$$W_{f,S}^Q = W_{f,S}^I \oplus Y_{f,S}^Q \oplus \bigoplus_{s \in S \cap \mathcal{K}(f)} W_{f,u_f \cup S \setminus s}^Q. \quad (17)$$

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

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

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

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



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

$$= Y_{f,S \cup k \setminus u_f}^Q \oplus \bar{W}_{f,S \cup k \setminus u_f} \oplus \bigoplus_{s \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} W_{f,S \cup k \setminus s}^Q \quad (18)$$

and  $W_{f,S}^I$

$$= W_{f,S}^Q \oplus Y_{f,S \cup k \setminus u_f}^I \oplus W_{f,S \cup k \setminus u_f}^I \oplus \bigoplus_{s \in \{\mathcal{S} \setminus u_f\} \cap \mathcal{K}(f)} \bar{W}_{f,S \cup k \setminus s} \quad (19)$$

for all sets  $S \in \bar{\mathcal{T}}_k(r)$  with  $u_f \in S$ . Now user  $U_k$  can mimic the user leader  $u_f$  and obtain the remaining subfiles, i.e.  $W_{f,S}$  for all sets  $S \in \bar{\mathcal{T}}_k(r)$  satisfying  $u_f \notin S$  from (16) and (17).

### B. Extension to arbitrary requested files

In previous subsection, we required all files to be requested. In this section, we show that the same scheme can be applied to any request if configured for  $K + 1$  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, we add the  $K + 1$ th “virtual” user to the system i.e.  $\mathcal{K} = \{1, \dots, K + 1\}$  and assume that the  $K + 1$ th user requests any of the files i.e.  $\mathbf{d}(K + 1) \in \mathcal{F}$ . Then, we apply the delivery strategy as described in previous subsections. In that case, every user obtains their desired file, with a rate

$$R = N \left( 1 - \frac{r}{K+1} \right). \quad (20)$$

Instead, if some files are not requested, the additional user requests the binary sum of all files not requested, i.e.  $W_0 \triangleq \bigoplus_{f \in \mathcal{F} \setminus \mathcal{N}_e(\mathbf{d})} W_f$  where  $\mathcal{N}_e(\mathbf{d})$  denotes the set of files requested. Now, we can replace the set of files in the system  $\mathcal{F}$  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 subfiles not requested as  $W_{0,S} \triangleq \bigoplus_{f \in \mathcal{F} \setminus \mathcal{N}_e(\mathbf{d})} W_{f,S}$ . Similarly, users can consider the cached subfiles in (2), as if the set of file was  $\mathcal{F}_0(\mathbf{d})$ , since

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

and, obtain the coded subfiles in (3) associated to file  $W_0$  as

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

Now, since we have replaced the set of files  $\mathcal{F}$  by  $\mathcal{F}_0(\mathbf{d})$ , and all files in  $\mathcal{F}_0(\mathbf{d})$  are requested, we can apply the delivery strategy as described in previous subsection, with a rate

$$R = (N_e(\mathbf{d}) + 1) \left( 1 - \frac{r}{K+1} \right). \quad (21)$$

Combining (20) and (21), we can write the rate as presented in Theorem 1.

## V. 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.

## 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] —, “Decentralized coded caching attains order-optimal memory-rate tradeoff,” *IEEE/ACM Transactions on Networking*, vol. 23, no. 4, pp. 1029–1040, Aug 2015.
- [3] 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.
- [4] 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.
- [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, June 2016.
- [6] 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.
- [7] H. Ghasemi and A. Ramamoorthy, “Improved lower bounds for coded caching,” in *IEEE ISIT*, June 2015, pp. 1696–1700.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] 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.
- [12] 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>
- [13] 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.
- [14] 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.
- [15] C. Y. Wang, S. S. Bidokhti, and M. Wigger, “Improved converses and gap-results for coded caching,” pp. 2428–2432, June 2017.
- [16] J. Gómez-Vilardebó, “Fundamental Limits of Caching: Improved Bounds with Coded Prefetching,” *ArXiv e-prints*, Dec. 2016.
- [17] —, “A novel centralized coded caching scheme with coded prefetching,” Dec. 2017. [Online]. Available: <http://www.cttc.es/publication/coded-caching/>