

# Multi-access coded caching: gains beyond cache-redundancy

Berksan Serbetci, Emanuele Parrinello and Petros Elia  
Communication Systems Department, EURECOM, Sophia Antipolis, France  
Email: {serbetci,parrinel,elia}@eurecom.fr

**Abstract**— The work considers the  $K$ -user cache-aided shared-link broadcast channel where each user has access to exactly  $z$  caches of normalized size  $\gamma$ , and where each cache assists exactly  $z$  users. For this setting, for two opposing memory regimes, we propose novel caching and coded delivery schemes which maximize the local caching gain, and achieve a *coding gain* larger than  $1 + K\gamma$  (users served at a time) despite the fact that the total cache redundancy remains  $K\gamma$  irrespective of  $z$ . Interestingly, when  $z = \frac{K-1}{K\gamma}$ , the derived optimal coding gain is  $K\gamma z + 1$ , matching the performance of a hypothetical scenario where each user has its own dedicated cache of size  $z\gamma$ .

## I. INTRODUCTION

A significant step toward understanding the fundamental limits of interference-limited caching networks was achieved in the work of Maddah-Ali and Niesen in [1] where an information theoretic study of the cache-aided broadcast channel (BC) revealed large improvements over the conventional caching approaches, as a result of a carefully designed cache placement phase which allowed a simultaneous delivery to many users at a time. In particular, the work in [1] considered a setting of a single server having access to a library of  $N$  unit-sized files and connected via a shared-link unit-capacity BC to  $K$  receiving users, where each user has access to its own dedicated cache of size  $M \leq N$ . In this setting, [1] showed that any possible set of requests by the  $K$  users, can be served with delay  $T = \frac{K(1-\gamma)}{K\gamma+1}$ , where  $\gamma = M/N$  is the normalized cache size. This revealed an ability to serve  $K\gamma + 1$  users at a time, where this number is commonly referred to as the *coding gain* or simply the *degrees-of-freedom (DoF)*. For this setting, the gain was proven in [2], [3] to be optimal under the constraint of uncoded cache placement. The coded caching ideas sparked considerable interest, resulting in a variety of related works that include [5], [6], [7], [8], [9], [10].

### A. Coded caching with multiple-access

Most of the works on coded caching consider scenarios where each user has its own dedicated cache. However in a variety of settings, such as different cellular networks, users can conceivably connect to multiple caches whose coverage areas may overlap. This consideration motivated the work in [11] which considered a similar  $K$ -user shared-link BC, where

This work was supported by the European Research Council under the EU Horizon 2020 research and innovation program / ERC grant agreement no. 725929.

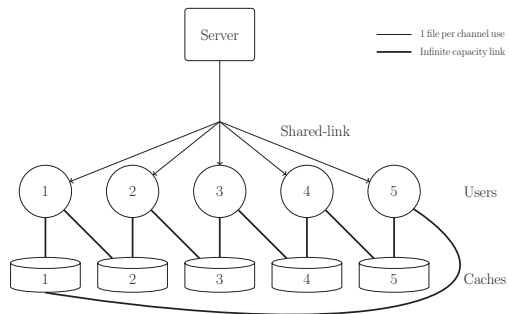


Fig. 1. Setting for  $K = 5$  and  $z = 2$ .

each user is assisted by exactly  $z > 1$  caches, and where each cache can serve exactly  $z$  users. In this context, the work in [11] provided a caching and delivery strategy whose centralized variant<sup>1</sup> achieves a worst-case delivery time of

$$T = \frac{K(1-z\gamma)}{K\gamma+1} \quad (1)$$

reflecting an ability to increase the local caching gain to  $z\gamma$  (meaning that each user sees a fraction  $z\gamma$  of each file), as well as an ability to preserve the coding gain to  $K\gamma + 1$ , by virtue of a scheme that served  $K\gamma + 1$  users at a time. For this same setting, in [12] the authors provide explicit designs for  $K = 4$ ,  $N = 4$ ,  $z = \{2, 3\}$ ,  $M = 1$  and  $K = 6$ ,  $N = 6$ ,  $z = 3$ ,  $M = 1$ , for which matching lower bounds are also developed to prove optimality of the schemes under uncoded cache placement. Finally, the authors provide a scheme for the extreme case of  $z = K - 1$ .

For this multi-access setting, we will show that the coding gain  $\frac{K(1-z\gamma)}{T}$  can exceed  $K\gamma + 1$ . In particular, for two opposing memory regimes, we propose two novel schemes which can serve, on average, more than  $K\gamma + 1$  users at a time. For the special case of  $z = \frac{K-1}{K\gamma}$ , the achieved gain is proven to be optimal under uncoded cache placement.

## II. SYSTEM MODEL AND PROBLEM DEFINITION

We consider a network where  $K$  users are connected via an error-free shared link to a server storing  $N$  ( $N \geq K$ ) files  $W^1, W^2, \dots, W^N$ . Each user has access to  $z$  out of  $K$  helper

<sup>1</sup>The work in [11] proposed a decentralized (stochastic) cache placement scheme, whose performance is slightly reduced over the easy-to-extend-to centralized variant whose delivery time we recorded above.

$$X_1 = \sum_{k=1}^K \left[ \sum_{j=1}^{\min\{k-z-1, z-1\}} \left[ \left\lfloor \frac{K - (k+z+j)}{2} \right\rfloor + 1 \right]_+ + \sum_{j=\max\{1, k-z\}}^{z-1} \left[ \left\lfloor \frac{K - 2z - 2j}{2} \right\rfloor + 1 \right]_+ \right],$$

$$X_2 = K \left[ \left\lfloor \frac{K-2-z}{3} \right\rfloor - z + 1 \right]_+ \left( K - 4z + 3 - 3 \left[ \left\lfloor \frac{K-2-z}{3} \right\rfloor - z + 1 \right]_+ \right), \quad S = \frac{K(K-2z+2)}{4}. \quad (2)$$

caches<sup>2</sup>, each of size  $M = N\gamma$  (units of file), where  $\gamma \in \{\frac{1}{K}, \frac{2}{K}, \dots, 1\}$ . We will use  $\mathcal{Z}_k$  to denote the content in cache  $k$ , and we will assume that each user has an unlimited capacity link to the caches it is connected to. Reflecting the assumption that each user has access to exactly  $z$  caches and each cache connects to exactly  $z$  users, we will consider, without loss of generality, the symmetric topology where each user  $k \in [K] \triangleq \{1, 2, \dots, K\}$  is associated to caches

$$\mathcal{C}_k \triangleq \langle k, k+1, \dots, k+z-1 \rangle \subseteq [K]$$

where in the above we use the notation  $\langle \mathcal{M} \rangle = \{m \mid m \text{ for } m \leq K; m-K \text{ for } m > K, m \in \mathbb{Z}^+, \forall m \in \mathcal{M}\}$ . A pictorial representation of the studied setting can be found in Figure 1.

The system works in two phases: the cache placement phase and the delivery phase. The first consists of filling the caches with the content of the library without knowledge of the users' demands. In the delivery phase, each user  $k$  requests a file from the library. We denote the index of such file by  $d_k$  and we collect all the indices of the requested files in the demand vector  $\mathbf{d} = (d_1, d_2, \dots, d_K)$ . The work focuses on the worst case where all users request different files. Upon reception of the demand vector  $\mathbf{d}$ , the server will transmit a message  $X$  that spans  $T$  units of time. Each user  $k$  will use  $X$  and their own available cache content  $\cup_{i \in \mathcal{C}_k} \mathcal{Z}_i$  to decode the desired file  $W^{d_k}$ . Our objective is to provide a caching and delivery scheme that reduces the delivery delay  $T$ .

### III. MAIN RESULTS

This section presents the main results. The proof of the following theorem follows from the scheme described in Section IV.

**Theorem 1.** *In the coded caching setting where each user is connected to  $z$  caches, when  $K\gamma = 2$ , the delivery time*

$$T = \frac{X_1 + X_2}{S} \quad (3)$$

*is achievable*<sup>3</sup>, where  $X_1, X_2$  and  $S$  are given in (2).

<sup>2</sup>Notice that the assumption of having as many users as caches can be relaxed to a more general case where the system has more users than caches, with a potential non uniform distribution of the users among the caches. Under these circumstances, the schemes described here can be combined with the non-uniformity adaptive coded caching scheme presented in [7] for the single transmitter setting.

<sup>3</sup>The above expression holds for the case where  $S(K-2z) - 4X_1$  is non negative and divisible by 3. This assumption can be removed at a small cost of increased scheme subpacketization.

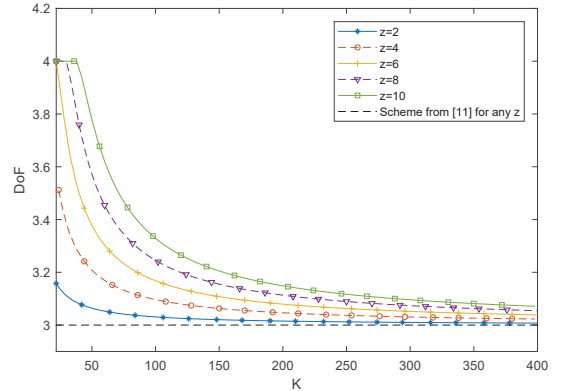


Fig. 2. DoF as a function of  $K$  for several  $z$ . Case of  $K\gamma = 2$ .

*Remark 1.* The scheme consists of transmissions serving either  $K\gamma + 1 = 3$  users or 4 users at a time. Figure 2 shows, for several values of  $z$ , the achievable DoF as a function of  $K$ . We can see how the DoF always exceeds  $K\gamma + 1 = 3$  and can approach 4 for some values of  $K$  and  $z$ .

We now characterize the optimal worst-case delivery time (under the assumption of uncoded cache placement) for the case of  $z = \frac{K-1}{K\gamma}$ , for any  $K\gamma$ .

**Theorem 2.** *In the addressed  $K$ -user caching network with access to an integer number  $z = \frac{K-1}{K\gamma}$  of caches of normalized size  $\gamma$ , the optimal delivery time, under the assumption of uncoded cache placement, takes the form*

$$T^* = \frac{K(1-\gamma z)}{K\gamma z + 1} = \frac{1}{K} \quad (4)$$

*corresponding to a DoF of  $K\gamma z + 1$  users served at a time.*

The scheme that achieves the above performance is described in Section V. The optimality — under uncoded cache placement — follows directly from the fact that the achieved performance matches the optimal performance (cf. [2],[3]) of a dedicated-cache coded caching setting (identical to that in [1]), where each cache has an augmented size equal to  $z\gamma$ .

### IV. CACHING AND DELIVERY SCHEME FOR $K\gamma = 2$

In this section we present a caching and delivery scheme for the case of  $K\gamma = 2$ . The scheme preserves the full local caching gain as in [11], and achieves a coding gain that strictly exceeds  $K\gamma + 1$ .

### A. Cache placement algorithm

In the cache placement phase, we first split each file  $W^n$  into  $S = \frac{K(K-2z+2)}{4}$  subfiles  $W_{\mathcal{T}}^n$  for each pair  $\mathcal{T} \triangleq \{\mathcal{T}_1, \mathcal{T}_2\}$  from the set

$$\Psi \triangleq \{\mathcal{T} : \mathcal{T}_1 \in [K-z], \mathcal{T}_2 \in [\mathcal{T}_1+z : 2 : \min\{K-z+\mathcal{T}_1, K\}]\} \quad (5)$$

where in the above we used the notation  $[a : b : c]$  to denote an ordered set of integers<sup>4</sup>, from  $a$  to  $c$ , in additive steps of  $b$ . After splitting the files, each cache  $k$  is filled as follows

$$\mathcal{Z}_k = \{W_{\mathcal{T}}^n \mid \forall n \in [N], \forall \mathcal{T} \ni k\}. \quad (6)$$

*a) Verifying the memory constraint:* To show that the cache placement satisfies the per-cache memory  $M = \frac{2N}{K}$ , we focus without loss of generality on cache 1, and note that the number of subfiles (per file) in this cache is  $\frac{K-2z}{2} + 1$ . Recalling that each such subfile is of size  $1/S$ , yields

$$\frac{N \left( \frac{K-2z}{2} + 1 \right)}{\frac{K(K-2z+2)}{4}} = \frac{2N}{K} = M$$

thus proving that the memory constraint is satisfied.

### B. Delivery scheme

The delivery scheme has two phases, where the first phase transmits XORs composed of 4 subfiles, while the second phase transmits XORs composed of  $K\gamma + 1 = 3$  subfiles.

*1) Phase 1:* Recall from above that any subfile  $W_{\mathcal{T}}^{d_k}$ ,  $k \in \mathcal{T} \in \Psi$  ( $\mathcal{T}_1 \in [k : 1 : \langle k+z-1 \rangle]$ ), is already available at one of the caches seen by the requesting user  $k \in [K]$ . For each user  $k \in [K]$ , the aim of this first phase is to serve the subfiles in the set

$$\left\{ W_{\mathcal{T}_1, \mathcal{T}_2}^{d_k} \mid \forall \{\mathcal{T}_1, \mathcal{T}_2\} \in \Psi : \mathcal{T}_1 \in \langle [k-z+1 : 1 : k-1] \rangle \cup \langle [k+z : 1 : k+2z-2] \rangle \right\}. \quad (7)$$

For any  $k \in [K]$ , let us define the following two sets

$$\begin{aligned} \Omega_{k,1} &\triangleq [k+1 : 1 : k+z-1] \\ \Omega_{k,2} &\triangleq [k-z+1 : 1 : k-1] \end{aligned} \quad (8)$$

and the set

$$B_{k,j} \triangleq [\Omega_{k,1}(j) + z : 2 : U_{k,j}] \quad (9)$$

where

$$U_{k,j} \triangleq \begin{cases} \langle \Omega_{k,2}(j) - z \rangle, & \text{if } \Omega_{k,2}(j) - z < 0 \\ K, & \text{otherwise.} \end{cases} \quad (10)$$

and where in the above we used the notation  $\Gamma(j)$  to denote the  $j$ -th element of an ordered set  $\Gamma$ . Next, for any  $k \in [K]$  and any  $j \in [z-1]$  we form the following XOR

$$\begin{aligned} X(k, j, m) &= W_{\langle \Omega_{k,1}(j) \rangle, B_{k,j}(m)}^{d_{\langle k-z+1 \rangle}} \oplus W_{\langle \Omega_{k,2}(j) \rangle, B_{k,j}(m)}^{d_k} \\ &\oplus W_{\langle \Omega_{B_{k,j}(m), 1}(j) \rangle, k}^{d_{\langle B_{k,j}(m) - z + 1 \rangle}} \oplus W_{\langle \Omega_{B_{k,j}(m), 2}(j) \rangle, k}^{d_{B_{k,j}(m)}}. \end{aligned} \quad (11)$$

<sup>4</sup>Note that  $b$  may be less than zero.

Creating the above XORs for every  $m \in [B_{k,j}]$  and every  $k \in [K]$ , spans the entire set of requested files in (7), and what we show below is that each component subfile (in the XORs) can be successfully decoded by its corresponding user.

*a) Decoding:* Consider any XOR as in (11) and let us focus on the subfiles  $W_{\Omega_{k,1}(j), B_{k,j}(m)}^{d_{\langle k-z+1 \rangle}}$  and  $W_{\Omega_{k,2}(j), B_{k,j}(m)}^{d_k}$  which are desired by users  $\langle k-z+1 \rangle$  and  $k$ , respectively. By the cache placement phase, we notice that the subfiles  $W_{\mathcal{T}_{B_{k,j}(m), z, 1}(j), k}^{d_{\langle B_{k,j}(m) - z + 1 \rangle}}$  and  $W_{\mathcal{T}_{b_{j,k}, z(m), z, 2}(j), k}^{d_{B_{k,j}(m)}}$  are both cached in cache  $k$ , thus enabling both users  $\langle k-z+1 \rangle$  and  $k$  to subtract these subfiles from  $X(k, j, m)$ . Next, we also notice that user  $\langle k-z+1 \rangle$  can cache out  $W_{\Omega_{k,2}(j), B_{k,j}(m)}^{d_k}$  since  $\Omega_{k,2}(j) \in \cup_{i \in \mathcal{C}_{\langle k-z+1 \rangle}} \mathcal{Z}_i$ . Similarly, user  $k$  can remove  $W_{\Omega_{k,1}(j), B_{k,j}(m)}^{d_{\langle k-z+1 \rangle}}$  from  $X(k, j, m)$  because  $\Omega_{k,1}(j) \in \cup_{i \in \mathcal{C}_k} \mathcal{Z}_i$ . Hence, we conclude that any XOR in (11) is decodable by both users  $k$  and  $\langle k-z+1 \rangle$ . In the same way, it can be shown that users  $\langle B_{k,j}(m) - z + 1 \rangle$  and  $B_{k,j}(m)$  can successfully decode their own requested subfiles.

*2) Phase 2:* We start by defining the set

$$\delta = \left[ z : 2 : \left\lfloor \frac{K-2-z}{3} \right\rfloor \right]$$

as well as the following set of triplets

$$\begin{aligned} \Theta = \left\{ (\theta_1, \theta_2, \theta_3) \mid \theta_1 = \delta(j), \theta_2 = 2\delta(j) + z + 2(i-1), \right. \\ \left. \theta_3 = \delta(j) + 2(i-1), j \in [\delta], i \in \left[ \frac{K-3\delta(j)-z}{2} \right] \right\}. \end{aligned} \quad (12)$$

For each triplet  $\theta \in \Theta$  and for each  $p \in [K]$ , we generate the following two XORs

$$\begin{aligned} Y_p(\theta, 1) &= W_{\langle \theta_2 - \theta_1 + p - 1 \rangle, \langle \theta_2 + p - 1 \rangle}^{d_p} \oplus W_{p, \langle \theta_2 + p \rangle}^{d_{\langle \theta_2 - \theta_1 - z + p \rangle}} \\ &\oplus W_{\langle z - 2 + p \rangle, \langle z - 2 + \theta_3 + p \rangle}^{d_{\langle \theta_2 + p - 1 \rangle}} \end{aligned} \quad (13)$$

$$\begin{aligned} Y_p(\theta, 2) &= W_{\langle \theta_2 - \theta_1 + p \rangle, \langle \theta_2 + p \rangle}^{d_p} \oplus W_{p, \langle \theta_2 + p \rangle}^{d_{\langle \theta_2 - \theta_1 - z + 1 + p \rangle}} \\ &\oplus W_{\langle z + p - 1 \rangle, \langle z + \theta_3 + p - 1 \rangle}^{d_{\langle \theta_2 + p - 1 \rangle}}. \end{aligned} \quad (14)$$

It can be shown (this is not done here due to lack of space) that, due to the structure of the XORs, as we go over all  $\theta \in \Theta, p \in [K]$ , no subfile is ever repeated. This allows us to conclude that the two phases successfully include all desired subfiles by all the users. As we did for phase 1, below we show that each subfile in the above XORs from (13) and (14) can be decoded successfully by their requesting user.

*a) Decoding:* For any  $p \in [K]$ , we will prove that the subfiles in  $Y_p(\theta, 1)$  can be decoded by their intended users. The decodability proof for  $Y_p(\theta, 2)$  will then follow directly.

User  $p$  can cache out subfiles  $W_{p, \langle \theta_2 + p \rangle}^{d_{\langle \theta_2 - \theta_1 - z + p \rangle}}$  and  $W_{\langle z - 2 + p \rangle, \langle z - 2 + \theta_3 + p \rangle}^{d_{\langle \theta_2 + p - 1 \rangle}}$  from  $Y_p(\theta, 1)$  since their subscripts  $p \in \cup_{i \in \mathcal{C}_p} \mathcal{Z}_i$  and  $\langle z - 2 + p \rangle \in \cup_{i \in \mathcal{C}_p} \mathcal{Z}_i$  correspond to the caches that user  $p$  is connected to. Next, we notice that user

$\langle \theta_2 - \theta_1 - z + p \rangle$  is connected to cache  $\langle \theta_2 - \theta_1 + p - 1 \rangle$  and thus it can cache out subfile  $W_{(\theta_2 - \theta_1 + p - 1), \langle \theta_2 + p - 1 \rangle}^{d_p}$ . The same user  $\langle \theta_2 - \theta_1 - z + p \rangle = \langle \delta(j) + p + 2(i - 1) \rangle$  has access to subfiles with subscripts in the set  $\langle [\delta(j) + p + 2(i - 1) : 1 : \delta(j) + p + 2(i - 1) + z - 1] \rangle$ . A relabelling of  $\langle \theta_3 + p + z - 2 \rangle$  to  $\langle \delta(j) + z + p - 2 + 2(i - 1) \rangle$  highlights that user  $\langle \theta_2 - \theta_1 - z + p \rangle$  can also remove  $W_{\langle z - 2 + p \rangle, \langle z - 2 + \theta_3 + p \rangle}^{d_{\langle \theta_2 + p - 1 \rangle}}$  from  $Y_p(\theta, 1)$ , and hence obtains its desired subfile  $W_{p, \langle \theta_2 + p \rangle}^{d_{\langle \theta_2 - \theta_1 - z + p \rangle}}$  successfully. Finally, we recall that user  $\langle \theta_2 + p - 1 \rangle$  has access to caches  $\mathcal{C}_{\langle \theta_2 + p - 1 \rangle} = \langle [\theta_2 + p - 1 : 1 : \theta_2 + p - 2 + z] \rangle$  and hence can successfully decode its desired subfile since it can cache out subfiles  $W_{(\theta_2 - \theta_1 + p - 1), \langle \theta_2 + p - 1 \rangle}^{d_p}$  and  $W_{p, \langle \theta_2 + p \rangle}^{d_{\langle \theta_2 - \theta_1 - z + p \rangle}}$ .

3) *Performance of the algorithm:* We observe that each generated XOR serves a different set of 3 or 4 subfiles, which can all be decoded. It can also be shown that the proposed delivery scheme successfully satisfies any demand vector  $\mathbf{d}$ . The proof is easy and is here omitted due to lack of space. Following the construction, we can readily count the total number of XORs transmitted during phase 1 and phase 2 to respectively be  $X_1$  and  $X_2$  from (2), which concludes the proof of the achievable delay in Theorem 1.

### C. Example

In this subsection we offer an example that may help to better understand the scheme. We consider the setting with parameters  $K = 10$ ,  $K\gamma = 2$ , and  $z = 2$ . For the sake of simplicity, we will use 0 to represent the index 10 and also, when describing a double index  $i, j$ , we will omit the comma.

In the placement phase, we first split each file, according to (5), into  $S = 20$  equally-sized subfiles with indices

$$\Psi = \{13, 15, 17, 19, 24, 26, 28, 20, 35, 37, 39, 46, 48, 40, 57, 59, 68, 60, 79, 80\}$$

and we then fill the caches according to (6) as follows

$$\begin{aligned} \mathcal{Z}_1 &= \{W_{13}^n, W_{15}^n, W_{17}^n, W_{19}^n, \forall n \in [N]\} \\ \mathcal{Z}_2 &= \{W_{24}^n, W_{26}^n, W_{28}^n, W_{20}^n, \forall n \in [N]\} \\ &\vdots \\ \mathcal{Z}_9 &= \{W_{19}^n, W_{39}^n, W_{59}^n, W_{79}^n, \forall n \in [N]\} \\ \mathcal{Z}_0 &= \{W_{20}^n, W_{40}^n, W_{60}^n, W_{80}^n, \forall n \in [N]\}. \end{aligned}$$

We notice that the cache placement guarantees an empty intersection  $\mathcal{Z}_k \cap \mathcal{Z}_{(k+1)} = \emptyset$  of any  $z = 2$  neighboring caches, and thus a full local caching gain ( $z\gamma = 0.4$ ).

In the delivery phase we consider the worst-case demand vector  $\mathbf{d} = (1, 2, \dots, 9, 0)$ . We will list the XORs of phase 1 and phase 2, but before doing that, let us offer some intuition on the design of the XORs of the first phase.

Let us consider a pair of users (say, users 0 and 1) that “see” a common cache (in this case, cache 1). For these two users we will create a generic XOR

$$W_{\sigma_1, \sigma_2}^0 \oplus W_{\tilde{\sigma}_1, \tilde{\sigma}_2}^1 \quad (15)$$

which will be combined with another XOR

$$W_{\tau_1, \tau_2}^3 \oplus W_{\tilde{\tau}_1, \tilde{\tau}_2}^4 \quad (16)$$

which is meant for another pair of users, say 3 and 4, that again share a common cache (cache 4). Combining the two XORs yields a new XOR  $X = W_{\sigma_1, \sigma_2}^0 \oplus W_{\tilde{\sigma}_1, \tilde{\sigma}_2}^1 \oplus W_{\tau_1, \tau_2}^3 \oplus W_{\tilde{\tau}_1, \tilde{\tau}_2}^4$  of 4 subfiles. To guarantee decoding for all, we will set  $\sigma_1 = \tilde{\sigma}_1 = 4$  to let user 3 and user 4 “cache out” from  $X$  the subfiles in (15) and similarly we set  $\tau_1 = \tilde{\tau}_1 = 1$  in order to let users 0 and 1 cache out the XOR in (16). Next, we choose  $\sigma_2 = 2$  so that user 1 can remove subfile  $W_{4,2}^0$  from  $X$  and  $\tilde{\sigma}_2 = 0$  to let user 0 remove subfile  $W_{4,0}^1$  from  $X$ . A similar choice of  $\tau_2$  and  $\tilde{\tau}_2$  will result in<sup>5</sup>

$$X = W_{24}^0 \oplus W_{40}^1 \oplus W_{15}^3 \oplus W_{13}^4. \quad (17)$$

The list of XORs sent during phase 1 is given below.

$$\begin{aligned} X(1, 1, 1) &= W_{24}^0 \oplus W_{40}^1 \oplus W_{15}^3 \oplus W_{13}^4 \\ X(1, 1, 2) &= W_{26}^0 \oplus W_{60}^1 \oplus W_{17}^5 \oplus W_{15}^6 \\ X(1, 1, 3) &= W_{28}^0 \oplus W_{80}^1 \oplus W_{19}^7 \oplus W_{17}^8 \\ X(2, 1, 1) &= W_{35}^1 \oplus W_{25}^2 \oplus W_{26}^4 \oplus W_{24}^5 \\ X(2, 1, 2) &= W_{37}^1 \oplus W_{27}^2 \oplus W_{28}^6 \oplus W_{26}^7 \\ X(2, 1, 3) &= W_{39}^1 \oplus W_{29}^2 \oplus W_{20}^8 \oplus W_{28}^9 \\ X(3, 1, 1) &= W_{46}^2 \oplus W_{26}^3 \oplus W_{37}^5 \oplus W_{35}^6 \\ X(3, 1, 2) &= W_{48}^2 \oplus W_{28}^3 \oplus W_{39}^7 \oplus W_{37}^8 \\ X(3, 1, 3) &= W_{40}^2 \oplus W_{20}^3 \oplus W_{13}^9 \oplus W_{39}^0 \\ X(4, 1, 1) &= W_{57}^3 \oplus W_{47}^4 \oplus W_{48}^6 \oplus W_{46}^7 \\ X(4, 1, 2) &= W_{59}^3 \oplus W_{39}^4 \oplus W_{40}^8 \oplus W_{48}^9 \\ X(5, 1, 1) &= W_{68}^4 \oplus W_{48}^5 \oplus W_{59}^7 \oplus W_{57}^8 \\ X(5, 1, 2) &= W_{60}^4 \oplus W_{40}^5 \oplus W_{15}^9 \oplus W_{59}^0 \\ X(6, 1, 1) &= W_{79}^5 \oplus W_{59}^6 \oplus W_{60}^8 \oplus W_{68}^9 \\ X(7, 1, 1) &= W_{80}^6 \oplus W_{60}^7 \oplus W_{17}^9 \oplus W_{79}^0. \end{aligned}$$

In phase 2, the  $X_2 = 20$  transmissions from (13) and (14) are cyclically generated as shown below.

$$\begin{aligned} Y_1(\theta, 1) &= W_{46}^1 \oplus W_{17}^3 \oplus W_{13}^6, \\ Y_2(\theta, 1) &= W_{57}^2 \oplus W_{28}^4 \oplus W_{24}^7 \\ Y_3(\theta, 1) &= W_{68}^3 \oplus W_{39}^5 \oplus W_{35}^8 \\ &\vdots \\ Y_0(\theta, 1) &= W_{35}^0 \oplus W_{60}^2 \oplus W_{20}^5 \\ \text{and} \\ Y_1(\theta, 2) &= W_{57}^1 \oplus W_{17}^4 \oplus W_{24}^6 \\ Y_2(\theta, 2) &= W_{68}^2 \oplus W_{28}^5 \oplus W_{35}^7 \\ Y_3(\theta, 2) &= W_{79}^3 \oplus W_{39}^6 \oplus W_{46}^8 \\ &\vdots \\ Y_0(\theta, 2) &= W_{46}^0 \oplus W_{60}^3 \oplus W_{13}^5. \end{aligned}$$

<sup>5</sup>This intuition can be generalized for  $z > 2$  and used as a baseline in the general description of the scheme presented in Section IV-B1.

In the end, we have  $S = 20$ ,  $X_1 = 15$  and  $X_2 = 20$  which gives

$$T = \frac{X_1 + X_2}{S} = \frac{35}{20}$$

and a coding gain  $K(1 - z\gamma)/T = 3.43$ .

## V. CODED CACHING FOR $K = K\gamma z + 1$

Corresponding to Theorem 2, we now present the optimal caching and delivery scheme for  $z = \frac{K-1}{K\gamma}$  for any  $K\gamma$ .

### A. Cache placement algorithm

In the cache placement phase, each file  $W^n$  is first split into  $K$  subfiles  $W_\phi^n$ ,  $\phi \in \Phi$  where each  $K\gamma$ -tuple  $\phi \triangleq \{\phi_1, \phi_2, \dots, \phi_{K\gamma}\}$  is drawn from the set

$$\Phi \triangleq \{\phi : \phi_1 \in [K], \phi_j = \langle \phi_{j-1} + z \rangle, \forall j \in [2 : K\gamma]\} \quad (18)$$

of size  $K$ . Then each cache  $k$  is filled as follows

$$\mathcal{Z}_k = \{W_\phi^n \mid \forall n \in [N], \forall \phi \ni k\} \quad (19)$$

forcing each integer  $k \in [K]$  to appear in  $\Phi$  exactly  $K\gamma$  times, thus guaranteeing that each cache stores exactly  $K\gamma$  subfiles from each file, thus respecting the cache-size constraint.

What the above placement also guarantees is that, by construction of the set  $\Phi$ , all subfiles of each file stored in any  $z$  consecutive caches, are different. This is due to the fact that any two elements of each  $K\gamma$ -tuple  $\phi \in \Phi$  have distance at least  $z$ . Thus, each user  $k$  has access to  $K\gamma z = K \frac{K-1}{Kz} z = K-1$  different subfiles of its requested file  $W^{d_k}$ . We denote by  $W_{\rho_k}^{d_k}$  the one remaining subfile desired by user  $k$ , for a specific  $K\gamma$ -tuple  $\rho_k = \Phi \setminus \{\cup_{i \in \mathcal{C}_k} \mathcal{Z}_i\}$ .

### B. Delivery and decoding

Upon reception of the demand vector  $\mathbf{d}$ , the server multicasts a single XOR

$$X = \bigoplus_{k \in [K]} W_{\rho_k}^{d_k} \quad (20)$$

to all users of the network. By virtue of the fact that each user is only missing a single subfile, we can deduce that each user  $k$  can cache out from  $X$  all  $K-1$  subfiles  $\{W_{\rho_j}^{d_j}\}_{j \in [K] \setminus \{k\}}$  to successfully decode its own requested subfile  $W_{\rho_k}^{d_k}$ . Consequently the total delivery time is naturally equal to  $T = |X| = \frac{1}{K}$ .

### C. Example

Let us consider the case of  $K\gamma = 3$ ,  $z = 2$  and  $K = 7$ . We first split each file into 7 equally sized subfiles with indices

$$\Phi = \{135, 136, 146, 246, 247, 257, 357\}.$$

and fill each cache, according to (19), as follows

$$\begin{aligned} \mathcal{Z}_1 &= \{W_{135}^n, W_{136}^n, W_{146}^n, \forall n \in [N]\} \\ \mathcal{Z}_2 &= \{W_{246}^n, W_{247}^n, W_{257}^n, \forall n \in [N]\} \end{aligned}$$

$$\mathcal{Z}_3 = \{W_{135}^n, W_{136}^n, W_{357}^n, \forall n \in [N]\}$$

$$\mathcal{Z}_4 = \{W_{146}^n, W_{246}^n, W_{247}^n, \forall n \in [N]\}$$

$$\mathcal{Z}_5 = \{W_{135}^n, W_{257}^n, W_{357}^n, \forall n \in [N]\}$$

$$\mathcal{Z}_6 = \{W_{136}^n, W_{146}^n, W_{246}^n, \forall n \in [N]\}$$

$$\mathcal{Z}_7 = \{W_{247}^n, W_{257}^n, W_{357}^n, \forall n \in [N]\}.$$

In the delivery phase we consider the delivery vector  $\mathbf{d} = (1, 2, \dots, 7)$ . We notice that the placement and topology jointly guarantee that each user is missing only a single subfile. For example, user 1 is only missing subfile  $W_{357}^1$ . Placing all these missing subfiles together, the server sends

$$X = W_{357}^1 \oplus W_{146}^2 \oplus W_{257}^3 \oplus W_{136}^4 \oplus W_{247}^5 \oplus W_{135}^6 \oplus W_{246}^7 \quad (21)$$

which guarantees that each user can cache out exactly 6 elements to decode their own subfile. The delay is  $T = \frac{1}{7}$  and the DoF of  $K\gamma z + 1 = 7$ .

## VI. CONCLUSION

For the multiple access coded caching problem, we have proposed novel coded caching schemes that achieve a coding gain that exceeds  $K\gamma + 1$ . To the best of our knowledge, in the context of worst-case delivery time and for  $N \geq K$ , this is the first instance of a fully-connected shared-link problem that experiences a gain larger than  $K\gamma + 1$ . An interesting extension of this work is its generalization to all possible values of  $K\gamma$ .

## REFERENCES

- [1] M. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. Inf. Theory*, May 2014.
- [2] K. Wan, D. Tuninetti, and P. Piantanida, "On the optimality of uncoded cache placement," in *Information Theory Workshop (ITW)*, IEEE, 2016.
- [3] Q. Yu, M. A. Maddah-Ali, and A. S. Avestimehr, "The exact rate-memory tradeoff for caching with uncoded prefetching," *IEEE Trans. Inf. Theory*, Feb 2017.
- [4] E. Lempiris, J. Zhang, and P. Elia, "Cache-aided cooperation with no CSIT," in *Proc. IEEE Int. Symp. on Inform. Theory (ISIT)*, June 2017.
- [5] T. X. Vu, S. Chatzinotas, and B. Ottersten, "Coded caching and storage planning in heterogeneous networks," in *2017 IEEE Wireless Communications and Networking Conference (WCNC)*, March 2017.
- [6] K. Wan, D. Tuninetti, M. Ji, and G. Caire, "Novel inter-file coded placement and d2d delivery for a cache-aided fog-ran architecture," 2018. [Online]. Available: <https://arxiv.org/pdf/1811.05498.pdf>
- [7] E. Parrinello, A. Ünsal, and P. Elia, "Fundamental limits of caching in heterogeneous networks with uncoded prefetching," *arXiv preprint https://arxiv.org/pdf/1811.06247.pdf*, 2018.
- [8] N. Zhang and M. Tao, "Fitness-aware coded multicasting for decentralized caching with finite file packetization," *IEEE Wireless Communications Letters*, Oct 2018.
- [9] H. Ghasemi and A. Ramamoorthy, "Algorithms for asynchronous coded caching," in *51st Asilomar Conference on Signals, Systems, and Computers*, Oct 2017.
- [10] E. Ozfatura, T. Rarris, D. Gndz, and O. Ercetin, "Delay-aware coded caching for mobile users," in *IEEE 29th Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, Sep. 2018.
- [11] J. Hachem, N. Karamchandani, and S. Diggavi, "Coded caching for multi-level popularity and access," *IEEE Trans. Inf. Theory*, May 2017.
- [12] K. S. Reddy and N. Karamchandani, "On the exact rate-memory trade-off for multi-access coded caching with uncoded placement," in *2018 International Conference on Signal Processing and Communications (SPCOM)*, July 2018.
- [13] E. Lempiris and P. Elia, "Achieving full multiplexing and unbounded caching gains with bounded feedback resources," in *Proc. IEEE Int. Symp. on Inform. Theory (ISIT)*, June 2018.
- [14] E. Lempiris and P. Elia, "Full coded caching gains for cache-less users," in *Information Theory Workshop (ITW)*, IEEE, 2018.