

Centralized Coded Caching with Heterogeneous Cache Sizes

Abdelrahman M. Ibrahim, Ahmed A. Zewail, and Aylin Yener

Wireless Communications and Networking Laboratory (WCAN)
 Electrical Engineering and Computer Science Department
 The Pennsylvania State University, University Park, PA 16802.
 {ami137,zewail}@psu.edu yener@engr.psu.edu

Abstract—Coded caching can improve fundamental limits of communication, utilizing storage memory at individual users. This paper considers a centralized coded caching system, introducing heterogeneous cache sizes at the users, i.e., the users’ cache memories are of different size. The goal is to design cache placement and delivery policies that minimize the worst-case delivery load on the server. To that end, the paper proposes an optimization framework for cache placement and delivery schemes which explicitly accounts for the heterogeneity of the cache sizes. We also characterize explicitly the optimal caching scheme, for the case where the sum of the users’ cache sizes is smaller than or equal to the library size.

I. INTRODUCTION

Caching [1]–[9] aims to alleviate network congestion during peak-traffic hours, by pushing data into the cache memories at the network edge during off-peak hours. The former is often called the delivery phase and the latter the placement phase. The contents requested by the users, are thus partially available at the users’ local cache memories and the remaining pieces need to be delivered, e.g., by a server. Reference [1] has introduced the fundamental concept of *coded caching*, where the placement and delivery phases are jointly optimized to minimize the delivery load, by creating multicast opportunities.

Extensive effort has transpired towards studying fundamental limits of coded caching in various setups since the introduction of coded caching systems with equal cache sizes in [1]. In particular, reference [2] has considered a decentralized system where the users populate their cache memories independently, i.e., the placement phase is decentralized. Despite lack of coordination, the server is able to create multicast transmissions and the resulting delivery load is close to the centralized delivery load in [1] for large networks. Several other network architectures with caching capabilities have also been considered, such as two-hop networks [5], device-to-device (D2D) networks [6], and interference networks [7], [8].

Today’s networks contain variety of devices at end users, e.g., laptops, smart phones, etc., with varying storage capabilities, which in turn motivates considering the realistic model where the end users have heterogeneous cache sizes. Reference [3] has extended the *decentralized* scheme in [2] to systems with heterogeneous cache sizes. Reference [4] has proposed a group-based delivery scheme that achieves lower delivery load for systems where the number of users is greater

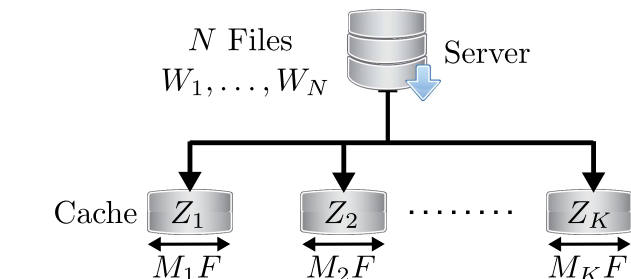


Fig. 1: Caching system with heterogeneous cache sizes.

than the number of files. For centralized caching however, the scheme in [1] cannot be directly extended to the setup with heterogeneous cache sizes. Thus, a new scheme is needed.

In this paper, we tackle this challenge and propose a centralized coded caching scheme in systems with heterogeneous cache sizes. We formulate and solve the optimization problem over the parameters of the cache placement and delivery schemes, in order to minimize the worst-case delivery load. We consider the set of feasible placement policies, \mathfrak{A} , under the uncoded prefetching assumption [9], i.e., the users cache uncoded pieces of the files. For given cache contents, we consider the set of delivery schemes, \mathfrak{D} , that includes all possible unicast/multicast transmissions, and guarantees that the users’ requests will be satisfied. For the case where library size is larger than or equal to accumulated users’ cache sizes, we derive the explicit caching scheme that yields the optimal worst-case delivery load under uncoded prefetching $R_{\mathfrak{A}}^*$. Finally, we compare the delivery load of our scheme with the lower bounds in [3], [4].

Notation: Vectors are represented by boldface letters, sets of policies are represented by calligraphic letters, e.g., \mathfrak{A} , \oplus refers to bitwise XOR operation, $|W|$ denotes cardinality of W , $\{\}$ denotes the empty set, $[K] := \{1, \dots, K\}$, and $2^{[K]}$ denotes the power set of $[K]$.

II. SYSTEM MODEL

We consider the caching system shown in Fig. 1, where a single server is connected to K users via a shared error-free multicast link [1]. A library $\{W_1, \dots, W_N\}$ of N files, each with size F bits, is stored at the server. We consider a

heterogeneous system, where user k is equipped with a cache memory of size $M_k F$ bits, i.e., M_k is the memory size of user k normalized by the file size. Without loss of generality, we assume that $M_1 \leq M_2 \leq \dots \leq M_K$. Additionally, we define $m_k = M_k/N$ to denote the memory size of user k normalized by the library size NF , i.e., $m_k \in [0, 1]$ for $M_k \in [0, N]$. The cache size vector is denoted by $\mathbf{M} = [M_1, \dots, M_K]$ and its normalized version by the library size is denoted by $\mathbf{m} = [m_1, \dots, m_K]$. We focus on the case, where the number of files is at least as many as the number of users, i.e., $N \geq K$.

The system operates over two phases: placement phase and delivery phase. In the placement phase, the server populates the users' cache memories without the knowledge of users' demands that will be made in the delivery phase. In particular, user k stores a subset Z_k of the files library, subject to its cache size constraint. In the delivery phase, user k requests a file W_{d_k} from the server. The users' demands are uniform and independent, i.e., the demand vector $\mathbf{d} = [d_1, \dots, d_K]$ consists of identical and independent uniform random variables over the files as in [1]. In order to serve the users' demands, the server transmits a sequence of unicast/multicast signals, $X_{\mathcal{T}, \mathbf{d}}$, to the users in $\mathcal{T} \in 2^{[K]} - \{\}$, which represent all possible transmission sets. Define the transmission variable $v_{\mathcal{T}} \in [0, 1]$ which represents the fraction of the file delivered to the users in \mathcal{T} via $X_{\mathcal{T}, \mathbf{d}}$. Thus, $X_{\mathcal{T}, \mathbf{d}}$ has length $v_{\mathcal{T}} F$ bits. At the end of the delivery phase, user k must be able to decode \hat{W}_{d_k} reliably. Formally, we have the following definition.

Definition 1. For a given normalized cache size vector \mathbf{m} , the delivery load $R(\mathbf{m})$ is said to be achievable if for every $\epsilon > 0$ and large enough F , there exists a caching scheme such that $\max_{\mathbf{d}, k \in [K]} \Pr(\hat{W}_{d_k} \neq W_{d_k}) \leq \epsilon$. Moreover, the infimum over all achievable delivery loads is denoted by $R^*(\mathbf{m})$. ■

In this work, we focus on uncoded prefetching where the users cache uncoded bits [9]. We also have that user k divides its cache size $M_k F$ equally over the files, since the demand is uniform, dedicating $m_k F$ bits to each file. We denote this class of cache placement policies by \mathfrak{A} . We denote the set of all delivery policies by \mathfrak{D} .

Definition 2. Under a placement policy in \mathfrak{A} , and any delivery policy in \mathfrak{D} , the worst-case delivery load is defined as $R_{\mathfrak{A}, \mathfrak{D}} := \max_{\mathbf{d}} R_{\mathbf{d}, \mathfrak{A}, \mathfrak{D}} = \sum_{\mathcal{T} \in 2^{[K]} - \{\}} v_{\mathcal{T}}$. Moreover, by taking the infimum over all caching schemes, we get $R_{\mathfrak{A}}^*$. ■

III. CACHING MODEL: THE THREE-USER CASE

For clarity of exposition, we first develop cache placement and delivery schemes for three users, i.e., $K = 3$. We then generalize to $K > 3$.

A. Placement Phase

Each file W_l is partitioned into 2^3 disjoint subfiles, which are labeled by the users storing them, as shown in Fig. 2. $\tilde{W}_{l, \mathcal{S}}$ denotes the subfile stored at the users in \mathcal{S} . For example, $\tilde{W}_{l, \{1,2\}}$ is stored at users 1 and 2, and $\tilde{W}_{l, \{\}}$ is not stored at any user. The set of all possible allocation sets is the power set of $\{1, 2, 3\}$. We assume that subfile sizes are the same for all

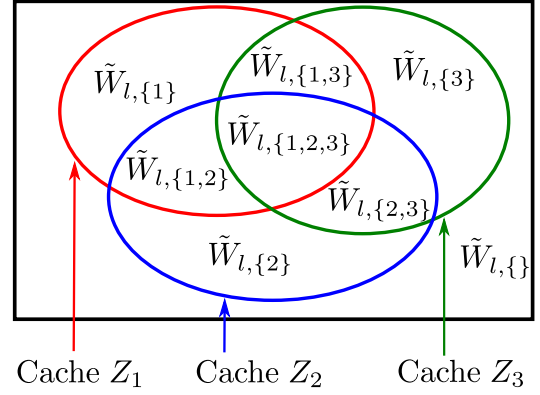


Fig. 2: The partitions of file W_l over the storage sets.

files, i.e., $|\tilde{W}_{l, \mathcal{S}}| = a_{\mathcal{S}} F$ bits, $\forall l \in [N]$, where the allocation variable $a_{\mathcal{S}} \in [0, 1]$ is the fraction of the file stored at \mathcal{S} . Hence, the cache content of user k is given by

$$Z_k = \bigcup_{l \in [N]} \bigcup_{\mathcal{S} \in 2^{[K]} : k \in \mathcal{S}} \tilde{W}_{l, \mathcal{S}}.$$

Moreover, the sets $\mathcal{S} \in 2^{[K]}$ define a partition of each file W_l , see Fig. 2. We have

$$a_{\{\}} + a_{\{1\}} + a_{\{2\}} + a_{\{3\}} + a_{\{1,2\}} + a_{\{1,3\}} + a_{\{2,3\}} + a_{\{1,2,3\}} = 1. \quad (1)$$

The fact that user k dedicates $m_k F$ bits of its cache memory to each file, implies that the cache sizes constraints can be expressed as follows

$$a_{\{1\}} + a_{\{1,2\}} + a_{\{1,3\}} + a_{\{1,2,3\}} \leq m_1, \quad (2)$$

$$a_{\{2\}} + a_{\{1,2\}} + a_{\{2,3\}} + a_{\{1,2,3\}} \leq m_2, \quad (3)$$

$$a_{\{3\}} + a_{\{1,3\}} + a_{\{2,3\}} + a_{\{1,2,3\}} \leq m_3. \quad (4)$$

Consequently, for $K = 3$, the set of feasible cache placement schemes is characterized by $a_{\mathcal{S}}, \mathcal{S} \in 2^{[K]}$ that satisfy (1)-(4).

B. Delivery Phase

Recall from Section II that $X_{\mathcal{T}, \mathbf{d}}$ is transmitted to the users in \mathcal{T} . At the end of transmission, user k must be able to reconstruct W_{d_k} from $X_{\mathcal{T}, \mathbf{d}}$, $k \in \mathcal{T}$ and its cache content Z_k .

1) *Structure of $X_{\mathcal{T}, \mathbf{d}}$:* The multicast signals are formed by XORing file pieces of equal size. For example, the multicast signal $X_{\{1,2\}, \mathbf{d}} = W_{d_1}^{\{1,2\}} \oplus W_{d_2}^{\{1,2\}}$, where $W_{d_j}^{\{1,2\}} \subset W_{d_j}$, and $|W_{d_1}^{\{1,2\}}| = |W_{d_2}^{\{1,2\}}| = v_{\{1,2\}} F$ bits. Similarly, we have

$$X_{\{1,3\}, \mathbf{d}} = W_{d_1}^{\{1,3\}} \oplus W_{d_3}^{\{1,3\}},$$

$$X_{\{2,3\}, \mathbf{d}} = W_{d_2}^{\{2,3\}} \oplus W_{d_3}^{\{2,3\}},$$

$$X_{\{1,2,3\}, \mathbf{d}} = W_{d_1}^{\{1,2,3\}} \oplus W_{d_2}^{\{1,2,3\}} \oplus W_{d_3}^{\{1,2,3\}}.$$

Additionally, a unicast transmission delivers the fraction of the requested file that is not available at the user's local cache and will not be delivered by the multicast transmissions. For example, the unicast transmission to user 1 is given by

$$X_{\{1\}, \mathbf{d}} = W_{d_1} - \bigcup_{\mathcal{S}: 1 \in \mathcal{S}} \tilde{W}_{d_1, \mathcal{S}} - W_{d_1}^{\{1,2\}} \bigcup W_{d_1}^{\{1,3\}} \bigcup W_{d_1}^{\{1,2,3\}},$$

where $\bigcup_{S:1 \in S} \tilde{W}_{d_1,S}$ is the fraction of the requested file W_{d_1} available locally at user 1 and $W_{d_1}^{\{1,2\}} \cup W_{d_1}^{\{1,3\}} \cup W_{d_1}^{\{1,2,3\}}$ is delivered to user 1 via multicast transmissions. Similarly, the remaining unicast signals are given by

$$X_{\{2\},d} = W_{d_2} - \bigcup_{S:2 \in S} \tilde{W}_{d_2,S} - W_{d_2}^{\{1,2\}} \cup W_{d_2}^{\{2,3\}} \cup W_{d_2}^{\{1,2,3\}},$$

$$X_{\{3\},d} = W_{d_3} - \bigcup_{S:3 \in S} \tilde{W}_{d_3,S} - W_{d_3}^{\{1,3\}} \cup W_{d_3}^{\{2,3\}} \cup W_{d_3}^{\{1,2,3\}}.$$

2) *Construction of $W_{d_j}^T$* : In order to ensure that each user in $\mathcal{T} - \{j\}$ is able to cancel $W_{d_j}^T$ from $X_{\mathcal{T},d}$, $W_{d_j}^T$ is constructed from the subfiles cached by all the users in $\mathcal{T} - \{j\}$, i.e., subfiles $\tilde{W}_{d_j,S}$ where $\mathcal{T} - \{j\} \subset S$ and $j \notin S$. In particular, $W_{d_j}^T$ is formed from the pieces $W_{d_j,S}^T$, i.e., the subsets of $W_{d_j}^T$ that are stored at the users in S . For example, the pieces composing $X_{\{1,2\}}$ are defined by

$$W_{d_1}^{\{1,2\}} = W_{d_1,\{2\}}^{\{1,2\}} \cup W_{d_1,\{2,3\}}^{\{1,2\}}, \quad W_{d_2}^{\{1,2\}} = W_{d_2,\{1\}}^{\{1,2\}} \cup W_{d_2,\{1,3\}}^{\{1,2\}}.$$

Similarly, the remaining pieces are defined by

$$W_{d_1}^{\{1,3\}} = W_{d_1,\{3\}}^{\{1,3\}} \cup W_{d_1,\{2,3\}}^{\{1,3\}}, \quad W_{d_1}^{\{1,2,3\}} = W_{d_1,\{2,3\}}^{\{1,2,3\}},$$

$$W_{d_3}^{\{1,3\}} = W_{d_3,\{1\}}^{\{1,3\}} \cup W_{d_3,\{1,2\}}^{\{1,3\}}, \quad W_{d_2}^{\{1,2,3\}} = W_{d_2,\{1,3\}}^{\{1,2,3\}},$$

$$W_{d_2}^{\{2,3\}} = W_{d_2,\{3\}}^{\{2,3\}} \cup W_{d_2,\{1,3\}}^{\{2,3\}}, \quad W_{d_3}^{\{1,2,3\}} = W_{d_3,\{1,2\}}^{\{1,2,3\}},$$

$$W_{d_3}^{\{2,3\}} = W_{d_3,\{2\}}^{\{2,3\}} \cup W_{d_3,\{1,2\}}^{\{2,3\}}.$$

Define the assignment variable $u_S^T \in [0, a_S]$ to represent the fraction of $\tilde{W}_{d_j,S}$ involved in $X_{\mathcal{T},d}$, i.e., $|W_{d_j,S}^T| = u_S^T F$ bits. Thus, each subfile $\tilde{W}_{d_j,S}$ is partitioned into $W_{d_j,S}^T$, and $W_{d_j,S}^{\{ \}}$ which denotes the remaining piece. We have

$$\tilde{W}_{d_3,\{1,2\}} = W_{d_3,\{1,2\}}^{\{1,3\}} \cup W_{d_3,\{1,2\}}^{\{2,3\}} \cup W_{d_3,\{1,2\}}^{\{1,2,3\}} \cup W_{d_3,\{1,2\}}^{\{ \}},$$

$$\tilde{W}_{d_2,\{1,3\}} = W_{d_2,\{1,3\}}^{\{1,2\}} \cup W_{d_2,\{1,3\}}^{\{2,3\}} \cup W_{d_2,\{1,3\}}^{\{1,2,3\}} \cup W_{d_2,\{1,3\}}^{\{ \}},$$

$$\tilde{W}_{d_1,\{2,3\}} = W_{d_1,\{2,3\}}^{\{1,2\}} \cup W_{d_1,\{2,3\}}^{\{2,3\}} \cup W_{d_1,\{2,3\}}^{\{1,2,3\}} \cup W_{d_1,\{2,3\}}^{\{ \}}.$$

3) *Delivery phase constraints*: Recall that we have defined the transmission variable, $v_{\mathcal{T}} \in [0, 1]$ as $|X_{\mathcal{T},d}|/F$ in Section II. We now describe the delivery scheme by constraints on the transmission variables $v_{\mathcal{T}}$, and the assignment variables u_S^T defined above in 2). First, in order to satisfy the users' demands, the signals sent by the server must *complete* the file requested by each user, which is represented by the following delivery constraints

$$v_{\{1\}} + v_{\{1,2\}} + v_{\{1,3\}} + v_{\{1,2,3\}} \geq 1 - m_1, \quad (5)$$

$$v_{\{2\}} + v_{\{1,2\}} + v_{\{2,3\}} + v_{\{1,2,3\}} \geq 1 - m_2, \quad (6)$$

$$v_{\{3\}} + v_{\{1,3\}} + v_{\{2,3\}} + v_{\{1,2,3\}} \geq 1 - m_3. \quad (7)$$

That is, m_k represents the local caching gain. Moreover, the structure of the transmitted pieces imposes the following constraints

$$v_{\{1,2\}} = u_{\{2\}}^{\{1,2\}} + u_{\{2,3\}}^{\{1,2\}} = u_{\{1\}}^{\{1,2\}} + u_{\{1,3\}}^{\{1,2\}}, \quad (8)$$

$$v_{\{1,3\}} = u_{\{3\}}^{\{1,3\}} + u_{\{2,3\}}^{\{1,3\}} = u_{\{1\}}^{\{1,3\}} + u_{\{1,2\}}^{\{1,3\}}, \quad (9)$$

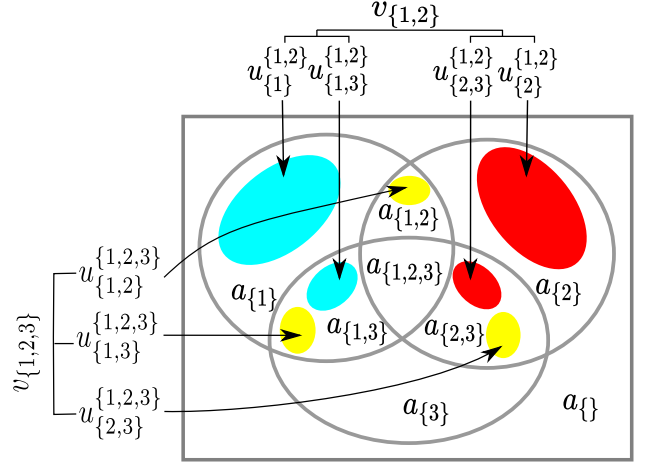


Fig. 3: The relation between transmission, allocation, and assignment variables.

$$v_{\{2,3\}} = u_{\{3\}}^{\{2,3\}} + u_{\{1,3\}}^{\{2,3\}} = u_{\{2\}}^{\{2,3\}} + u_{\{1,2\}}^{\{2,3\}}, \quad (10)$$

$$v_{\{1,2,3\}} = u_{\{2,3\}}^{\{1,2,3\}} = u_{\{1,3\}}^{\{1,2,3\}} = u_{\{1,2\}}^{\{1,2,3\}}. \quad (11)$$

The relationship between $v_{\mathcal{T}}$, the allocation variables, a_S , and the assignment variables, u_S^T , is further illustrated in Fig. 3. We observe that the transmission variable $v_{\{1,2\}}$ is limited by the amount of side information stored at user 1 and not available at user 2, i.e., $a_{\{1\}} + a_{\{1,3\}}$, and the amount of side information stored at user 2 and not available at user 1, i.e., $a_{\{2\}} + a_{\{2,3\}}$. Additionally, $v_{\{1,2\}}$ and $v_{\{1,2,3\}}$ are both limited by $a_{\{1,3\}}$ and $a_{\{2,3\}}$. Hence, we have the following side information constraints

$$v_{\{1,2\}} + v_{\{1,2,3\}} \leq a_{\{1\}} + a_{\{1,3\}}, \quad (12)$$

$$v_{\{1,2\}} + v_{\{1,2,3\}} \leq a_{\{2\}} + a_{\{2,3\}}. \quad (13)$$

Similarly, the other side information constraints are given by

$$v_{\{1,3\}} + v_{\{1,2,3\}} \leq a_{\{1\}} + a_{\{1,2\}}, \quad (14)$$

$$v_{\{1,3\}} + v_{\{1,2,3\}} \leq a_{\{3\}} + a_{\{2,3\}}, \quad (15)$$

$$v_{\{2,3\}} + v_{\{1,2,3\}} \leq a_{\{2\}} + a_{\{1,2\}}, \quad (16)$$

$$v_{\{2,3\}} + v_{\{1,2,3\}} \leq a_{\{3\}} + a_{\{1,3\}}. \quad (17)$$

Moreover, subfile $\tilde{W}_{d_3,\{1,2\}}$ can be utilized in the multicast transmissions that include user 3, i.e., $X_{\{1,3\},d}$, $X_{\{2,3\},d}$, $X_{\{1,2,3\},d}$. The following constraint prevents transmitting redundant bits to user 3

$$u_{\{1,2\}}^{\{1,3\}} + u_{\{1,2\}}^{\{2,3\}} + u_{\{1,2\}}^{\{1,2,3\}} \leq a_{\{1,2\}}. \quad (18)$$

Similarly, the constraints that prevent transmitting redundant bits to users 1 and 2 are given by

$$u_{\{1,3\}}^{\{1,2\}} + u_{\{1,3\}}^{\{2,3\}} + u_{\{1,3\}}^{\{1,2,3\}} \leq a_{\{1,3\}}, \quad (19)$$

$$u_{\{2,3\}}^{\{1,2\}} + u_{\{2,3\}}^{\{1,3\}} + u_{\{2,3\}}^{\{1,2,3\}} \leq a_{\{2,3\}}. \quad (20)$$

In summary, given the allocation variables, the set of feasible delivery schemes, i.e., schemes that satisfy all users' requests,

in a three-user system, is characterized by (5)-(20), $0 \leq u_S^T \leq a_S$, and $0 \leq v_T \leq 1$.

IV. CACHING MODEL: THE K -USER CASE

In this section, we generalize the cache placement and delivery schemes to systems with K users.

A. Placement phase

In the general case, each file W_l is partitioned into 2^K subfiles, and the set of feasible cache placement policies for a given \mathbf{m} , $\mathfrak{A}(\mathbf{m})$, is defined as

$$\left\{ \mathbf{a} \in [0, 1]^{2^K} \mid \sum_{S \in 2^{[K]}} a_S = 1, \sum_{S \in 2^{[K]} : k \in S} a_S \leq m_k, \forall k \in [K] \right\}, \quad (21)$$

where the allocation vector \mathbf{a} represents the collection of allocation variables a_S , $S \in 2^{[K]}$. The cache placement procedure is summarized in Algorithm 1, for a given allocation vector $\mathbf{a} \in \mathfrak{A}(\mathbf{m})$.

Algorithm 1 Cache placement procedure

Input: $\{W_1, \dots, W_N\}$ and \mathbf{a}

Output: $Z_k, k \in [K]$

- 1: **for** $l \in [N]$ **do**
 - 2: Partition file W_l into subfiles $\tilde{W}_{l,S}, S \in 2^{[K]}$ such that $|\tilde{W}_{l,S}| = a_S F$.
 - 3: **end for**
 - 4: **for** $k \in [K]$ **do**
 - 5: $Z_k \leftarrow \bigcup_{l \in [N]} \bigcup_{S \in 2^{[K]} : k \in S} \tilde{W}_{l,S}$
 - 6: **end for**
-

B. Delivery phase

In general, $X_{\mathcal{T}, \mathbf{a}} = \bigoplus_{j \in \mathcal{T}} W_{d_j}^T$, where $|W_{d_j}^T| = v_{\mathcal{T}} F$ bits, $\forall j \in \mathcal{T}$. Additionally, a multicast transmission to the users in \mathcal{T} is constrained by the side information stored at the sets

$$\mathcal{B}_j^T := \left\{ S \in 2^{[K]} : \mathcal{T} - \{j\} \subset S, j \notin S \right\}, \forall j \in \mathcal{T}, \quad (22)$$

where \mathcal{B}_j^T represents the sets containing the side information stored at $\mathcal{T} - \{j\}$ and not available at user j . The set of all storage sets related to \mathcal{T} is denoted by

$$\mathcal{B}^T := \bigcup_{j \in \mathcal{T}} \mathcal{B}_j^T. \quad (23)$$

Moreover, for $|\mathcal{T}| \geq 2$, each piece $W_{d_j}^T$ is partitioned over the sets $S \in \mathcal{B}_j^T$, i.e.,

$$W_{d_j}^T = \bigcup_{S \in \mathcal{B}_j^T} W_{d_j, S}^T, \quad (24)$$

where $|W_{d_j, S}^T| = u_S^T F$ bits.

For given \mathbf{m} and \mathbf{a} , the set of feasible delivery schemes is denoted by $\mathfrak{D}(\mathbf{m}, \mathbf{a})$, and characterized by the transmission and assignment variables that satisfy the following conditions

$$\sum_{S \in \mathcal{B}_j^T} u_S^T = v_{\mathcal{T}}, \forall \mathcal{T} \in 2^{[K]} - \{j\}, \forall j \in \mathcal{T}, \quad (25)$$

Algorithm 2 Delivery procedure

Input: $\mathbf{d}, \mathbf{a}, \mathbf{u}, \mathbf{v}$, and $\tilde{W}_{l,S}, S \in 2^{[K]}, l \in [N]$

Output: $X_{\mathcal{T}, \mathbf{a}}, \mathcal{T} \in 2^{[K]} - \{j\}$

- {Partitioning}
 - 1: **for** $\{S \in 2^{[K]} : 1 \leq |S| \leq K - 1\}$ **do**
 - 2: **for** $\{j \in [K] : j \notin S\}$ **do**
 - 3: Partition $\tilde{W}_{d_j, S}$ into $W_{d_j, S}^T, \{\mathcal{T} \in 2^{[K]} : j \in \mathcal{T}, \mathcal{T} \cap S \neq \{j\}\}$, such that $|W_{d_j, S}^T| = u_S^T F$ and $W_{d_j, S}^{\{j\}}$ is the remaining segment.
 - 4: **end for**
 - 5: **end for**
 - {Delivery scheme}
 - 6: **for** $\mathcal{T} \in 2^{[K]} - \{j\}$ **do**
 - 7: **if** $\mathcal{T} = \{j\}$ **then**
 - 8: $W_{d_j}^{\{j\}} := W_{d_j} - \bigcup_{S: j \in S} \tilde{W}_{d_j, S} - \bigcup_{\mathcal{T}'} \bigcup_S W_{d_j, S}^{\mathcal{T}'}$
 - 9: $X_{\{j\}, \mathbf{a}} \leftarrow W_{d_j}^{\{j\}}$ {Unicast transmissions}
 - 10: **else**
 - 11: $W_{d_j}^T := \bigcup_{S \in \mathcal{B}_j^T} W_{d_j, S}^T$
 - 12: $X_{\mathcal{T}, \mathbf{a}} \leftarrow \bigoplus_{j \in \mathcal{T}} W_{d_j}^T$ {Multicast transmissions}
 - 13: **end if**
 - 14: **end for**
-

$$\sum_{\mathcal{T} \in 2^{[K]} - \{j\} : k \in \mathcal{T}} v_{\mathcal{T}} \geq 1 - m_k, \forall k \in [K], \quad (26)$$

$$\sum_{\mathcal{T} \in 2^{[K]} - \{j\} : \{j\} \cup S' \subset \mathcal{T}} v_{\mathcal{T}} \leq \sum_{S \in 2^{[K]} : S' \subset S, j \notin S} a_S, \forall j \notin S', \quad (27)$$

$$\sum_{\mathcal{T} \in 2^{[K]} - \{j\} : j \in \mathcal{T}, \mathcal{T} \cap S \neq \{j\}, |\mathcal{T}| \leq |S| + 1} u_S^T \leq a_S, \forall j \notin S, \quad (28)$$

$$0 \leq u_S^T \leq a_S, \forall \mathcal{T} \in 2^{[K]} - \{j\}, \forall S \in \mathcal{B}^T, \quad (29)$$

$$0 \leq v_{\mathcal{T}} \leq 1, \forall \mathcal{T} \in 2^{[K]} - \{j\}. \quad (30)$$

The delivery scheme is summarized in Algorithm 2, for given allocation (\mathbf{a}), assignment (\mathbf{u}), and transmission vectors (\mathbf{v}).

V. OPTIMIZATION AND PERFORMANCE

In Sections III-A and IV-A, we have illustrated that a cache placement policy is completely characterized by the allocation vector \mathbf{a} , which represents the fraction of files stored exclusively at each subset of users $S \subset [K]$. Sections III-B and IV-B have demonstrated that a delivery scheme is characterized by the pair (\mathbf{u}, \mathbf{v}) , where the size and structure of the transmitted signals are represented by the transmission vector \mathbf{v} and the assignment vector \mathbf{u} , respectively. Next, we formulate an optimization problem to minimize the worst-case delivery load, i.e., the sum of transmission variables, by jointly optimizing the cache placement and delivery schemes.

A. Optimizing the Caching Scheme

For a given normalized memory vector \mathbf{m} , the following optimization problem characterizes the minimum worst-case

delivery load $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m})$ and the optimal caching scheme, i.e., the optimal values for \mathbf{a} , \mathbf{v} , and \mathbf{u} .

$$\text{OI: } \min_{\mathbf{a}, \mathbf{u}, \mathbf{v}} \sum_{\mathcal{T} \in 2^{[K]} - \{\emptyset\}} v_{\mathcal{T}} \quad (31a)$$

$$\text{subject to } \mathbf{a} \in \mathfrak{A}(\mathbf{m}), \quad (31b)$$

$$(\mathbf{u}, \mathbf{v}) \in \mathfrak{D}(\mathbf{m}, \mathbf{a}), \quad (31c)$$

where $\mathfrak{A}(\mathbf{m})$ is the set of feasible allocation vectors defined in (21) and $\mathfrak{D}(\mathbf{m}, \mathbf{a})$ is the set of feasible assignment and transmission vectors defined by (25)-(30). Note that $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m}) \geq R_{\mathfrak{A}}^*(\mathbf{m}) \geq R^*(\mathbf{m})$, where R^* , $R_{\mathfrak{A}}^*$ are defined in Section II.

Remark 1. In the case of equal cache sizes, $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m})$ is equal to the worst-case delivery of the MaddahAli-Niesen caching scheme in [1], which was shown to be optimal for uncoded prefetching in [9]. Moreover, the solution obtained from the optimization problem is equivalent to the memory sharing solution proposed in [1]. ■

In order to evaluate the performance of the proposed caching scheme, we compare $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m})$ to the following lower bound on $R^*(\mathbf{m})$ which is derived in [4]. For given K , N , \mathbf{M} , and $M_1 \leq \dots \leq M_K$, R^* is lower bounded by

$$\max_{s \in [K], l \in [\lceil \frac{N}{s} \rceil]} \frac{N - (N - Kl)^+}{l} - \frac{s \sum_{i=1}^{s+\gamma} M_i + \gamma(N - ls)^+}{l(s + \gamma)}, \quad (32)$$

where $\gamma := \min\left\{\left(\lceil \frac{N}{l} \rceil - s\right)^+, K - s\right\}$, and $(x)^+ := \max\{0, x\}$. This bound is tighter than the cut-set bound [3], given by

$$\max_{s \in \{\min\{K, N\}\}} \left\{ s - \frac{N \sum_{i=1}^s m_i}{\lfloor N/s \rfloor} \right\}. \quad (33)$$

B. Special Case: $\sum_{i=1}^K m_i \leq 1$

For the case where the sum of the cache sizes is less than or equal to the library size, i.e., $\sum_{i=1}^K m_i \leq 1$, we characterize explicitly the optimal solution of (31) as follows.

Proposition 1. For $\sum_{i=1}^K m_i \leq 1$ and $m_1 \leq \dots \leq m_K$, the minimum worst-case delivery load is given by

$$R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m}) = K - \sum_{j=1}^K (K - j + 1)m_j,$$

where $a_{\{j\}}^* = m_j$, $v_{\{j\}}^* = 1 - \sum_{i=1}^{j-1} m_i - (K - j + 1)m_j$, and $v_{\{i,j\}}^* = u_{\{i\}}^* = u_{\{j\}}^* = \min\{a_{\{i\}}^*, a_{\{j\}}^*\}$. ■

Proof. (Outline) The sum of the constraints in (26) gives

$$\begin{aligned} R_{\mathfrak{A},\mathfrak{D}}(\mathbf{m}) &\geq K - \sum_{i=1}^K m_i - \sum_{\mathcal{T} \in 2^{[K]} - \{\emptyset\}; |\mathcal{T}| \geq 2} (|\mathcal{T}| - 1) v_{\mathcal{T}}, \\ &\geq K - \sum_{j=1}^K (K - j + 1)m_j, \end{aligned}$$

since $\sum_{\mathcal{T} \in 2^{[K]} - \{\emptyset\}; |\mathcal{T}| \geq 2} (|\mathcal{T}| - 1) v_{\mathcal{T}} \leq \sum_{j=1}^K (K - j) m_j$,

for $m_1 \leq \dots \leq m_K$. Moreover, for $\sum_{i=1}^K m_i \leq 1$, $a_{\mathcal{S}}^*$, $v_{\mathcal{T}}^*$, and $u_{\mathcal{S}}^*$, represent a feasible solution that achieves the lower bound, i.e., $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m}) = K - \sum_{j=1}^K (K - j + 1)m_j$. ■

Using the genie-aided approach in [9], we can show that $R_{\mathfrak{A}}^*(\mathbf{m}) = R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m})$, which means that the caching scheme in Proposition 1 is optimal under uncoded prefetching.

Proposition 2. For $\sum_{i=1}^K m_i \leq 1$ and $m_1 \leq \dots \leq m_K$,

$$R_{\mathfrak{A}}^*(\mathbf{m}) \geq K - \sum_{j=1}^K (K - j + 1)m_j, \quad (34)$$

which implies $R_{\mathfrak{A}}^*(\mathbf{m}) = R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m})$. ■

Proof. (Outline) From the genie-aided approach in [9], we get

$$\begin{aligned} R_{\mathfrak{A}}^*(\mathbf{m}) &\geq \min_{\mathbf{a} \in \mathfrak{A}(\mathbf{m})} \sum_{j=1}^K \sum_{\mathcal{S} \in 2^{[K]}; \{1, \dots, j\} \cap \mathcal{S} = \{\emptyset\}} a_{\mathcal{S}}, \\ &= \min_{\mathbf{a} \in \mathfrak{A}(\mathbf{m})} K a_{\{\emptyset\}} + \sum_{j=1}^{K-1} j \sum_{\mathcal{S} \in 2^{[K]}; \{j+1, \dots, K\} \cap \mathcal{S} = \{\emptyset\}} a_{\mathcal{S}}, \end{aligned}$$

and (34) is obtained by solving the dual linear program. ■

C. Discussion

In our proposed scheme, we have formed $X_{\mathcal{T},\mathbf{a}}$ by using file pieces of equal size. Next, we remark that this restriction does not lead to a performance degradation.

Remark 2. A delivery scheme with $\tilde{X}_{\mathcal{T},\mathbf{a}} = \bigoplus_{j \in \mathcal{T}} \varphi_{v_{\mathcal{T}}} (W_{d_j}^T)$, where $\varphi_{v_{\mathcal{T}}}(\cdot)$ pads zeros to subfiles with size less than $v_{\mathcal{T}}F$ bits, is equivalent to a delivery scheme in \mathfrak{D} and both yield the same delivery load. For example, a multicast signal $\tilde{X}_{\{1,2\},\mathbf{a}} = \varphi_{v_{\{1,2\}}} (W_{d_1}^{\{1,2\}}) \oplus \varphi_{v_{\{1,2\}}} (W_{d_2}^{\{1,2\}})$, where $a_{\{2\}} > a_{\{1\}}$ and $\varphi_{v_{\{1,2\}}}(\cdot)$ appends $a_{\{2\}} - a_{\{1\}}$ zeros to $W_{d_2}^{\{1,2\}}$, is equivalent to a multicast signal $X_{\{1,2\},\mathbf{a}}$ and a unicast signal $X_{\{2\},\mathbf{a}}$, with sizes $a_{\{1\}}F$ bits, and $(a_{\{2\}} - a_{\{1\}})F$ bits, respectively. ■

Utilizing Remark 2 enables us to provide the following proposition comparing the performance of our propose scheme to that of the decentralized scheme in [3].

Proposition 3. For given N , K , \mathbf{m} and $m_1 \leq \dots \leq m_K$,

$$R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m}) \leq R_{dec}(\mathbf{m}) = \sum_{i=1}^K \prod_{j=1}^i (1 - m_j). \quad (35)$$

Proof. The decentralized placement scheme [3] is represented by $a_{\mathcal{S}} = \prod_{i \in \mathcal{S}} m_i \prod_{i \in \mathcal{S}^c} (1 - m_i)$, which belongs to \mathfrak{A} . The decentralized delivery scheme [3] is equivalent to a delivery scheme in \mathfrak{D} by Remark 2. Thus, the decentralized scheme in [3] is a feasible (but not necessarily optimal) solution to the optimization problem in (31). ■

VI. NUMERICAL RESULTS

In this section, we first present a numerical example that explains how to interpret the solution of the optimization

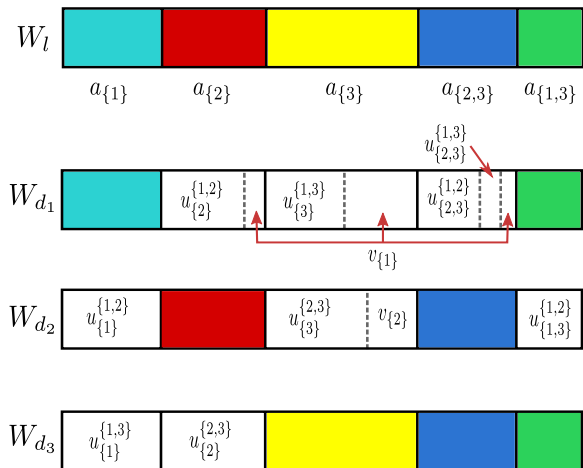


Fig. 4: Illustration of the optimal solution of (31) for a three-user system with cache sizes $m = [0.3, 0.4, 0.6]$.

problem in (31). Then, we compare numerically the centralized delivery load obtained from (31) with the decentralized delivery load and the existing lower bounds.

Example 1. Consider a caching system with $K = N = 3$, and $\mathbf{M} = [0.9, 1.2, 1.8]$, i.e., $\mathbf{m} = [0.3, 0.4, 0.6]$. The caching scheme obtained from (31), is characterized as follows.

Placement phase: Every file $W^{(l)}$ is divided into five subfiles, such that $a_{\{1\}} = 0.1958$, $a_{\{2\}} = 0.2042$, $a_{\{3\}} = 0.3$, $a_{\{2,3\}} = 0.1958$, and $a_{\{1,3\}} = 0.1042$, as shown in Fig. 4.

Delivery phase: We have the following transmissions

- 1) $X_{\{1,2\},a}$ delivers the pieces corresponding to $u_{\{2\}}^{\{1,2\}} = 0.17457$, $u_{\{2,3\}}^{\{1,2\}} = 0.12543$ to user 1, and $u_{\{1\}}^{\{1,2\}} = 0.1958$, $u_{\{1,3\}}^{\{1,2\}} = 0.1042$ to user 2, i.e., $v_{\{1,2\}} = 0.3$.
- 2) $X_{\{1,3\},a}$ delivers the pieces corresponding to $u_{\{3\}}^{\{1,3\}} = 0.15265$, $u_{\{2,3\}}^{\{1,3\}} = 0.04315$ to user 1, and $u_{\{1\}}^{\{1,3\}} = 0.1958$ to user 3, i.e., $v_{\{1,3\}} = 0.1958$.
- 3) $X_{\{2,3\},a}$ delivers the piece corresponding to $u_{\{3\}}^{\{2,3\}} = 0.2024$ to user 2, and the piece corresponding to $u_{\{2\}}^{\{2,3\}} = 0.2024$ to user 3, i.e., $v_{\{2,3\}} = 0.2024$.
- 4) The unicast transmissions to users 1 and 2 are represented by $v_{\{1\}} = 0.2024$ and $v_{\{2\}} = 0.0958$, respectively.

At the end of transmission, user k is able to decode file W_{d_k} , and the minimum worst-case delivery load $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m}) = 1$. The breakdown of the requested files is illustrated in Fig. 4. ■

In Fig. 5, we compare the centralized delivery load $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m})$, with the decentralized delivery load $R_{\text{dec}}(\mathbf{m})$, the lower bound in (32), the cut-set bound in (33), and the genie-aided bound for uncoded prefetching in (34), for $K = 4$, and $m_k = 0.9^{4-k} m_4$. We observe that $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m})$ coincides with the genie-aided lower bound for $\sum_{k=1}^4 m_k \leq 1$. Moreover, $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m})$ coincides with the cut-set bound for large memories, and $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m}) \leq R_{\text{dec}}(\mathbf{m})$.

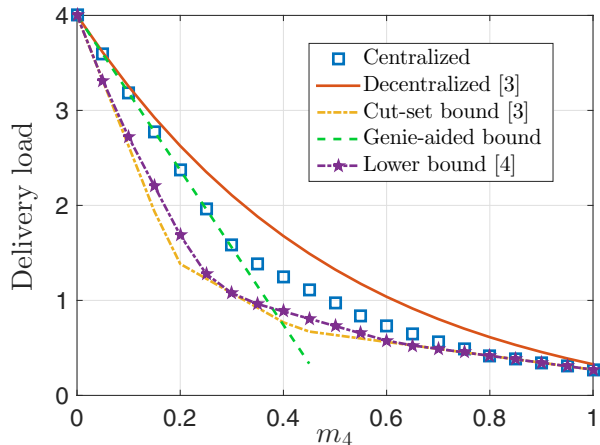


Fig. 5: Comparing $R_{\mathfrak{A},\mathfrak{D}}^*(\mathbf{m})$, $R_{\text{dec}}(\mathbf{m})$, cut-set and genie-aided lower bounds, for $K = 4$, and $m_k = 0.9 m_{k+1}$.

VII. CONCLUSIONS

In this paper, we have considered a centralized coded caching system with heterogeneous cache sizes. We have proposed a caching scheme taking into account the unequal cache sizes. We have put forward an optimization framework for minimizing the worst-case delivery load by optimizing the parameters of the caching scheme, for given cache memory sizes. We have shown that our caching scheme is optimal under uncoded prefetching for $\sum_{k=1}^K m_k \leq 1$. We have illustrated the caching scheme by a numerical example and compared our centralized delivery load with the decentralized delivery load derived in [3], and the lower bounds in [3], [4].

This work is a step towards developing caching schemes for systems with practical considerations, such as wireless systems where the users have distinct download rates and non-uniform resources. Future directions include optimizing over cache sizes and consideration of incomplete/inaccurate demand profiles.

REFERENCES

- [1] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. Info. Theory*, vol. 60, no. 5, pp. 2856–2867, Mar. 2014.
- [2] —, "Decentralized coded caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Trans. Netw.*, vol. 23, no. 4, pp. 1029–1040, Aug. 2015.
- [3] S. Wang, W. Li, X. Tian, and H. Liu, "Coded caching with heterogenous cache sizes," *arXiv:1504.01123*, 2015.
- [4] M. M. Amiri, Q. Yang, and D. Gündüz, "Decentralized coded caching with distinct cache capacities," *arXiv:1611.01579*, 2016.
- [5] N. Karamchandani, U. Niesen, M. A. Maddah-Ali, and S. N. Diggavi, "Hierarchical coded caching," *IEEE Trans. Info. Theory*, vol. 62, no. 6, pp. 3212–3229, Apr. 2016.
- [6] M. Ji, G. Caire, and A. F. Molisch, "Fundamental limits of caching in wireless D2D networks," *IEEE Trans. Info. Theory*, vol. 62, no. 2, pp. 849–869, Feb. 2016.
- [7] F. Xu, M. Tao, and K. Liu, "Fundamental tradeoff between storage and latency in cache-aided wireless interference networks," *arXiv:1605.00203*, 2016.
- [8] N. Naderializadeh, M. A. Maddah-Ali, and A. S. Avestimehr, "Fundamental limits of cache-aided interference management," *arXiv:1602.04207*, 2016.
- [9] Q. Yu, M. A. Maddah-Ali, and A. S. Avestimehr, "The exact rate-memory tradeoff for caching with uncoded prefetching," *arXiv:1609.07817*, 2016.