# A Multi-Antenna Coded Caching Scheme with Linear Subpacketization

MohammadJavad Salehi*, Antti Tölli*, Seyed Pooya Shariatpanahi†

*Center for Wireless Communications, University of Oulu, Oulu, Finland.

†School of Electrical and Computer Engineering, University of Tehran, Tehran, Iran.

{fist_name.last_name}@oulu.fi; p.shariatpanahi@ut.ac.ir

*Abstract*—Exponentially growing subpacketization is known to be a major issue for practical implementation of coded caching, specially in networks with multi-antenna communication setups. We provide a new coded caching scheme for such networks, which requires linear subpacketization and is applicable to any set of network parameters, as long as the multi-antenna gain $L$ is larger than or equal to the global caching gain $t$. Our scheme includes carefully designed cache placement and delivery algorithms; which are based on circular shift of two generator arrays in perpendicular directions. It also achieves the maximum possible degrees of freedom of $t + L$, during any transmission interval.

*Index Terms*—Coded Caching, Multi-Antenna Communications, Linear Subpacketization

## I. INTRODUCTION

Network data traffic has been subject to continuous growth during the past years. The total global IP (Internet Protocol) data volume is estimated to exceed $4.8$ Zettabytes ($10^{21}$ bytes) by 2022, from which 71 percent is expected to pass through wireless networks [1]. Different applications contribute to the wireless data traffic and each of them requires specific networking Key Performance Indicators (KPIs) such as data rate, delay and reliability. With the introduction of new application types for 5G and beyond (e.g. autonomous vehicles, immersive viewing and massive machine-type communications), extreme advancements for all KPI requirements are expected [2], [3]. This has imposed serious challenges in various network layers and solving them is one of the main recent research trends.

Among various networking KPIs, data rate is still of prominent importance. This is mainly due to video applications, as they are expected to account for $82\%$ of the global IP data traffic by 2022 [1]. 5G networks promote data rates of Gigabits per second; and further increase in the achievable data rate will still be a key driver in future wireless networks [3]. However, increasing wireless data rate is quite challenging and requires new resources to be used. 5G networks are introducing new frequency bands for cellular communications; from which mm-Wave bands are of much interest as they not only provide larger bandwidth, but also enable cell sizes to be decreased (resulting in better frequency reuse) and larger spatial gain of multi-antenna communications to be achieved [4], [5].

There is another important resource, recently proposed as a promising enabler of increased data rate for future networks.

This idea, originally proposed in [6], is known as Coded Caching and enables a global caching gain, proportional to the total cache size in the network, to be achieved in addition to the local caching gain at each cache location. This results in a new resource, i.e. storage, to become available for data networks; which is specially inspiring as the storage prices are constantly declining [7]. Interestingly, coded caching suits well for a majority of video-based applications, for which there is a prime request time (there are time intervals with high request rate) and uneven popularity distribution (a small set of files are requested repeatedly). Also as will be discussed later, it is shown that coded caching gain is additive with multi-antenna gain; making it even more desirable for future networks.

Despite its benefits, coded caching still requires major issues to be solved, before it can be practically implemented. In this paper we target one such issue, known as the large subpacketization requirement. The problem is that the number of smaller parts each file should be split into, known as subpacketization, grows exponentially with respect to the user count $K$; making coded caching implementation infeasible, even for moderate network sizes [8]. Considering networks with multi-antenna communication setups, we show that linear subpacketization growth is indeed possible, as long as the multi-antenna gain $L$ is larger than or equal to the coded caching gain $t$. Specifically, we introduce a novel coded caching scheme, which requires linear subpacketization with respect to primary network parameters $K$, $L$, $t$; to achieve the largest possible degrees of freedom of $t + L$ during any single transmission[1]. The sole feasibility condition of $L \geq t$ enables the scheme to be applied to a large class of networks; and is in line with the recent trend of using larger antenna arrays. This is a concrete solution to the subpacketization issue of coded caching schemes, making coded caching one step closer to practical implementation in next-generation networks.

In this paper, we use $[K]$ to denote $\{1, 2, ..., K\}$ and $[i : j]$ to represent $\{i, i + 1, ..., j\}$. Boldface upper- and lower-case letters denote matrices and vectors, respectively. $\mathbf{V}[i, j]$ refers to the element at the $i$-th row and $j$-th column of matrix $\mathbf{V}$. Sets are denoted by calligraphic letters. For two sets $\mathcal{A}$ and $\mathcal{B}$, $\mathcal{A} \backslash \mathcal{B}$ is the set of elements in $\mathcal{A}$ which are not in $\mathcal{B}$; and $|\mathcal{A}|$ represents the number of elements in $\mathcal{A}$.

[1]In this paper we assume $t$ does not scale with $K$. If $t$ scales with $K$, the growth in subpacketization will be quadratic.

## II. SYSTEM MODEL

We consider a multiple input, single output (MISO) broadcast setup, in which a single server communicates with $K$ users over a shared wireless link with the capacity of $f$ bits per channel use. The server has $L$ transmitting antennas and each user is equipped with a single antenna. Full channel state information (CSI) is available at the server; and it has access to a library of $N \geq K$ files, denoted by $\mathcal{F}$. Each file $W \in \mathcal{F}$ has a size of $f$ bits, and each user is equipped with a cache memory of size $Mf$ bits. For simplicity, we use a normalized data unit and drop $f$ in our subsequent notations.

The system operation consists of two distinct phases, placement and delivery. During the placement phase, which takes place at the low network traffic time, cache memories of the users are filled by data from the files in $\mathcal{F}$. This in done in accordance with a cache placement algorithm, which operates without any prior knowledge of file request probabilities in the delivery phase. We use $\mathcal{Z}(k)$ to denote the cache contents of user $k$, after the placement phase is completed.

At the beginning of the delivery phase, each user $k$ reveals its requested file $W(k) \in \mathcal{F}$. Let us define the demand set as $\mathcal{D} = \{W(k) \mid k \in [K]\}$. Based on $\mathcal{D}$ and in accordance with a delivery algorithm, the server builds $S$ transmission vectors $\mathbf{x}(1), \mathbf{x}(2), ..., \mathbf{x}(S)$, each with dimensions $L \times 1$ ($S$ is a design parameter depending on network parameters). Transmission vectors are then transmitted in a TDMA fashion, using the array of $L$ antennas. After $\mathbf{x}(s)$ is transmitted, user $k$ receives

$$y_k(s) = \mathbf{h}_k^T \mathbf{x}(s) + w_k(s) , \qquad (1)$$

where $\mathbf{h}_k \in \mathbb{C}^L$ denotes the $L \times 1$ channel vector (from $L$ transmitting antennas); and $w_k(s) \sim \mathbb{CN}(0,1)$ is the observed noise at user $k$ during transmission interval $s$. Transmission vectors are built such that each user $k$ can decode its requested file $W(k)$, using $\mathcal{Z}(k)$ (its locally cached data) together with $y_k(1), y_k(2), ..., y_k(S)$ (data received from the channel). Let us denote the set of users targeted by $\mathbf{x}(s)$ as $\mathcal{T}(s)$, for which we have $\mathcal{T}(s) \subseteq [K]$ and $|\mathcal{T}(s)| = t+L$. We assume zero-forcing beamformers $\mathbf{v}_{\mathcal{R}}$ are used to build $\mathbf{x}(s)$, where $\mathcal{R} \subseteq \mathcal{T}(s)$ and $|\mathcal{R}| = t+1$; and $\mathbf{v}_{\mathcal{R}}$ is built such that $\|\mathbf{v}_{\mathcal{R}}\| = 1$ and

$$\begin{aligned} \mathbf{h}_k^T \mathbf{v}_{\mathcal{R}} \neq 0 \qquad & k \in \mathcal{R} , \\ \mathbf{h}_k^T \mathbf{v}_{\mathcal{R}} = 0 \qquad & \mathcal{T}(S) \backslash \mathcal{R} . \end{aligned} \qquad (2)$$

We also assume that during downlink training, the server sends orthogonal demodulation pilots precoded by $\mathbf{v}_{\mathcal{R}}$; so that each user $k$ is able to estimate the equivalent channels $\mathbf{h}_k^T \mathbf{v}_{\mathcal{R}}, \forall \mathcal{R}$.

Delivery time $T$ is defined as the time required for all users to successfully decode their requested files. Cache placement and delivery algorithms should be designed such that the worst case delivery time (with respect to $\mathcal{D}$) is minimized. Let us denote the worst case delivery time by $T^*$. Following the common practice in the literature, we assume each user requests a different file, in order to find $T^*$. For simplicity, we also use the notation $A \equiv W(1)$, $B \equiv W(2)$, etc., in the examples provided in this paper.

As cache placement is done without any knowledge of file request probabilities, an efficient strategy is to store equal-sized data portions of all files in the cache memory of each user. Thereby, every user has $\frac{M}{N}$ of each file in its cache memory, and should receive the rest $(1 - \frac{M}{N})$ of its requested file from the server. This results in a total data size of $K(1 - \frac{M}{N})$ to be transmitted over the channel. Let us define the global cache ratio (coded caching gain) $t$ as the total cache size in the network normalized by the number of files, i.e. $t = \frac{KM}{N}$; and assume $t$ is an integer. Then the sum rate of the communication, denoted by $R^*$, is defined as

$$R^* = \frac{K(1 - \frac{t}{K})}{T^*} . \qquad (3)$$

As the channel capacity is one (normalized) data unit per channel use, the symmetric rate also represents how many users benefit from each transmission. So we use the term Degree of Freedom (DoF) equivalent to $R^*$. The goal is then to design cache placement and delivery algorithms such that DoF is maximized.

## III. STATE-OF-THE-ART

### A. Coded Caching

Coded caching is originally proposed by Maddah-Ali and Niesen in [6], where it is shown that DoF of $t+1$ is achievable with subpacketization $\binom{K}{t}$. This scheme is later extended in various directions; e.g. decentralized, hierarchical and multi-server coded caching [9]–[11]. Interestingly, in [11] it is shown that coded caching and multi-server gains are additive; and so DoF of $t + L$ is achievable with $L$ transmitting servers. However, the scheme of [11] requires larger subpacketization of

$$\binom{K}{t} \binom{K - t - 1}{L - 1} . \qquad (4)$$

Following the same concept, multi-antenna coded caching with zero-forcing beamformers is later introduced in [12], [13]. Optimized beamformers are then used in [14], to improve the performance at finite-SNR regime. In [15], interesting methods based on two design parameters $\alpha, \beta$ are introduced to reduce the optimized beamformer design complexity. However the subpacketization is further increased to

$$\frac{(\alpha - 1)!}{(\delta - 1)!(\beta - 1)!(t + \beta)!^{\delta - 1}} \binom{K}{t} \binom{K - t - 1}{L - 1} , \qquad (5)$$

where $\delta = \frac{t+\alpha}{t+\beta}$. In summary, the original scheme of [6] and its extensions for multi-antenna setups require exponentially growing subpacketization (with respect to $K$ and for fixed $\frac{M}{N}$), which makes the implementation infeasible even for moderate values of $K$ [8]. Consequently, reducing subpacketization without decreasing DoF has been studied in the literature, both for single- and multi-antenna coded caching.

### B. Subpacketization in Single-Antenna Coded Caching

Subpacketization is well-studied for single-antenna setups. In [16] it is shown that decentralized schemes need exponential subpacketization to achieve any sub-linear rate, for constant $\frac{M}{N}$ as $K \to \infty$. In [17] Placement Delivery Array (PDA) is presented as a systematic approach to reduce subpacketization in centralized schemes. It is shown that the original scheme of [6] is in fact a PDA-driven scheme; and is optimal among a symmetric class of schemes known as $g$-regular PDA.

Following [17], in [18] it is shown that for a constant rate $R^*$, a PDA resulting in linear subpacketization does not exists. In [19] a sub-exponential subpacketization scheme for fixed $R^*$ and $\frac{M}{N}$ is proposed. In [20] Ruzsa-Szemerédi graphs are used to design coded caching schemes with linear subpacketization as $K \to \infty$, but with non-constant $R^*$.

### C. Subpacketization in Multi-Antenna Coded Caching

Subpacketization is less studied for multi-antenna setups. Most notable work on this topic is [8], in which it is shown that if $\frac{K}{L}$ and $\frac{t}{L}$ are both integers, any single-antenna scheme with subpacketization $g(K,t)$ has a multi-antenna counterpart; with subpacketization $g(\frac{K}{L}, \frac{t}{L})$ and without any DoF loss ($g$ is a general function). For example, the scheme of [6] can be applied to multi-antenna setups, with subpacketization $\binom{K/L}{t/L}$. Unfortunately, the scheme of [8] suffers DoF loss (and also increased subpacketization), if either $\frac{K}{L}$ or $\frac{t}{L}$ is non-integer. Specifically, DoF is reduced by a multiplicative factor (gap), that is bounded above by 2 when $L > t$, and by $\frac{3}{2}$ when $L < t$.

In [21] it is shown that subpacketization can be traded-off with the performance; and a new approach is introduced for selecting subpacketization in a more flexible manner. The results are however limited to the specific case of $K = t+L$. In [22] joint reduction of CSI and subpacketization requirements is considered; and it shown that subpacketization of $L_c \binom{K_c}{t}$ is achievable, where $L_c = \frac{L+t}{t+1}$ and $K_c = \frac{K}{L_c}$. However, the proposed scheme requires both $L_c$ and $K_c$ to be integers; making it applicable to a very specific set of network parameters. Moreover, it results in a DoF loss by a factor of $(1 - \frac{t}{K})$.

### D. Our Contribution

We provide a coded caching scheme with DoF $t + L$ and subpacketization $K \times (t+L)$, for any network with $L \geq t$. The provided scheme requires linear subpacketization with respect to all network parameters $K$, $L$ and $t$, as long as the multi-antenna gain is larger than or equal to the coded caching gain.

## IV. CACHE PLACEMENT

Cache placement is based on placement matrices introduced in [21], which are also special cases of PDA [17]. A placement matrix $\mathbf{V}$ is a $P \times K$ binary matrix ($P$ can be any integer for which $\frac{Pt}{K}$ is an integer); with $\sum_p \mathbf{V}[p,k] = \frac{Pt}{K}, \forall k \in [K]$ and $\sum_k \mathbf{V}[p,k] = t, \forall p \in [P]$. Here we use a special placement matrix $\mathbf{V}$ with $P = K$, in which the first row has $t$ consecutive one elements (other elements are zero); and for the other rows, each row is a circular shift of the previous row by one unit. Given $\mathbf{V}$, we split each file $W$ into $P = K$ smaller parts $W_p$, and each part $W_p$ into $t + L$ smaller parts $W_p^q$. Then for every $p \in [K], k \in [K]$, if $\mathbf{V}[p,k] = 1$, $W_p^q$ is stored in the cache memory of user $k$, $\forall W \in \mathcal{F}, q \in [t + L]$.

**Example 1.** *Assume* $K = 6$, $t = 2$, $L = 3$. $\mathbf{V}$ *is built as*

$$\mathbf{V} = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \tag{6}$$

*and subpacketization is* $6 \times 5 = 30$. *Cache content of user 1 is*
$$\mathcal{Z}(1) = \{W_1^1, W_1^2, W_1^3, W_1^4, W_1^5,$$
$$W_6^1, W_6^2, W_6^3, W_6^4, W_6^5 \mid W \in \mathcal{F}\},$$
*and cache content of other users can be written accordingly.*

## V. DELIVERY

### A. Graphical Representation

Before formal description of the delivery algorithm, we provide a graphical illustration of its operation for the network of Example 1. The delivery algorithm operates in $K = 6$ rounds and at each round, $K - t = 4$ transmission vectors are built; resulting in $S = 24$ total transmission intervals. We also assume the demand set is $\mathcal{D} = \{A, B, C, D, E, F\}$, and ignore the modulation effect for notation clarity.

Graphical illustration for the first and second rounds are provided in Figures 1 and 2, respectively. In both figures, each matrix column represents a user and each row stands for a file part index. For example, the first column represents user one, and the first row stands for the first part of all files; i.e. $W_1^q$, $\forall W \in \mathcal{F}, q \in [t + L]$. A lightly shaded entry in the matrix means the data part is cached at the respective user. For example, $W_1^q, W_6^q$ are stored at user 1, $\forall W \in \mathcal{F}, q \in [t + L]$. Clearly, the cache placement indicated by Figures 1 and 2 is equivalent to the placement matrix $\mathbf{V}$ provided in Example 1.

Consequently, a dark shaded entry indicates which index of the requested file is sent to the respective user, during the given transmission interval. For example, in Figure 1a the entries $(3, 1), (3, 2), (1, 3), (1, 4), (1, 5)$ are dark shaded, which means the first transmission vector at round 1, i.e. $\mathbf{x}(1)$, includes $A_3^1, B_3^1, C_1^1, D_1^1, E_1^1$; and $\mathcal{T}(1) = [1 : 5]$. With $L = 3$ antennas, each part can be nulled out at two users; and we have

$$\mathbf{x}(1) = A_3^1 \mathbf{v}_{\{1,3,4\}} + B_3^1 \mathbf{v}_{\{2,3,4\}}$$
$$+ C_1^1 \mathbf{v}_{\{1,2,3\}} + D_1^1 \mathbf{v}_{\{1,2,4\}} + E_1^1 \mathbf{v}_{\{1,2,5\}}. \tag{7}$$

Note that all superscripts are set to 1, as no data is transmitted prior to $\mathbf{x}(1)$. According to (1) and (2), user 1 receives

$$y_1(1) = A_3^1 \mathbf{h}_1^T \mathbf{v}_{\{1,3,4\}} + \underline{C_1^1 \mathbf{h}_1^T \mathbf{v}_{\{1,2,3\}}}$$
$$+ \underline{D_1^1 \mathbf{h}_1^T \mathbf{v}_{\{1,2,4\}}} + \underline{E_1^1 \mathbf{h}_1^T \mathbf{v}_{\{1,2,5\}}} + w_1(1). \tag{8}$$

From Example 1, we know that $W_1^1 \in \mathcal{Z}(1)$, $\forall W \in \mathcal{F}$. Also, based to the system model, user 1 can estimate $\mathbf{h}_1^T \mathbf{v}_{\mathcal{R}}$, $\forall \mathcal{R}$. This means user 1 can reconstruct and remove the underlined terms from its received signal in (8); and finally decode $A_3^1$ interference-free. Similarly, users 2, 3, 4, 5 can decode $B_3^1$, $C_1^1$, $D_1^1$, $E_1^1$ respectively, resulting in DoF of $t + L = 5$ for the first transmission interval.

The next transmission vectors in round 1, i.e. $\mathbf{x}(2)$-$\mathbf{x}(4)$, are built by circular shift of $\mathbf{x}(1)$ elements over the non-shaded cells of the grid, in two perpendicular directions. Specifically, the first two terms of $\mathbf{x}(1)$ are shifted vertically, while the other three terms are shifted horizontally. This procedure is depicted in Figures 1b-1d. So $\mathbf{x}(2)$ is built as

$$\mathbf{x}(2) = A_4^1 \mathbf{v}_{\{1,4,5\}} + B_4^1 \mathbf{v}_{\{2,4,5\}}$$
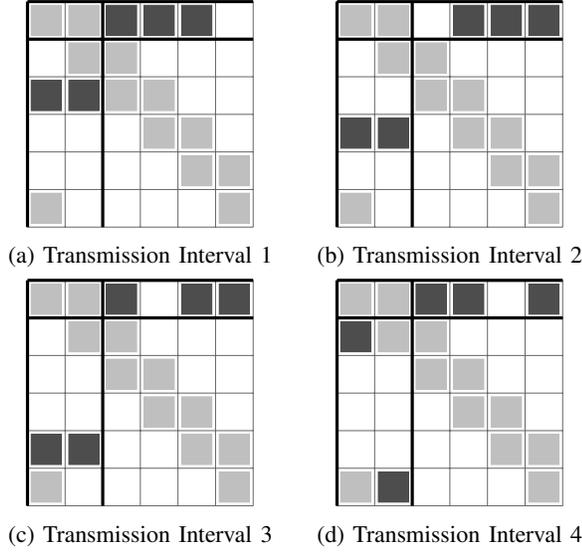$$+ D_1^2 \mathbf{v}_{\{1,2,4\}} + E_1^2 \mathbf{v}_{\{1,2,5\}} + F_1^1 \mathbf{v}_{\{1,2,6\}}, \tag{9}$$

(a) Transmission Interval 1     (b) Transmission Interval 2

(c) Transmission Interval 3     (d) Transmission Interval 4

Fig. 1: Graphical Illustration of the First Round



(a) Transmission Interval 5     (b) Transmission Interval 6

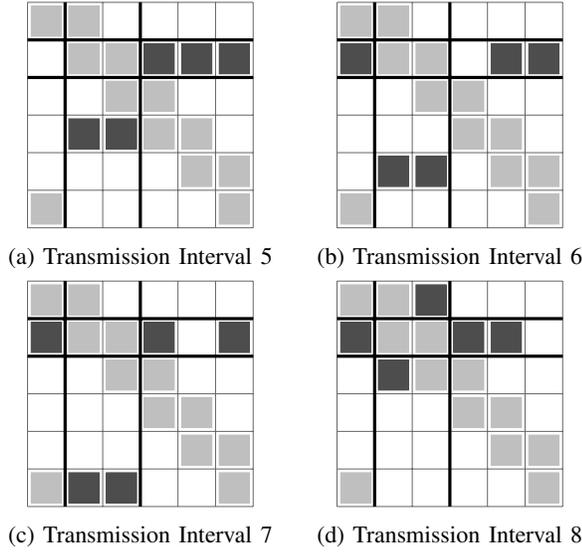(c) Transmission Interval 7     (d) Transmission Interval 8

Fig. 2: Graphical Illustration of the Second Round

where the superscripts for $D_1$ and $E_1$ are updated to 2, as $D_1^1$ and $E_1^1$ were transmitted by $\mathbf{x}(1)$. Similarly, $\mathbf{x}(3)$ and $\mathbf{x}(4)$ are built as

$$
\begin{aligned}
\mathbf{x}(3) =& A_5^1 \mathbf{v}_{\{1,5,6\}} + B_5^1 \mathbf{v}_{\{2,5,6\}} \\
&+ E_1^3 \mathbf{v}_{\{1,2,5\}} + F_1^2 \mathbf{v}_{\{1,2,6\}} + C_1^2 \mathbf{v}_{\{1,2,3\}} \;, \\
\mathbf{x}(4) =& A_2^1 \mathbf{v}_{\{1,2,3\}} + B_6^1 \mathbf{v}_{\{1,2,6\}} \\
&+ F_1^3 \mathbf{v}_{\{1,2,6\}} + C_1^3 \mathbf{v}_{\{1,2,3\}} + D_1^3 \mathbf{v}_{\{1,2,4\}} \;.
\end{aligned}
\tag{10}
$$

The second transmission round includes $\mathbf{x}(5)$-$\mathbf{x}(8)$, which are built by diagonal shift (simultaneous circular shift of one unit to the right and down) of $\mathbf{x}(1)$-$\mathbf{x}(4)$; as shown in Figure 2. Similarly, the third round is built by diagonal shift of $\mathbf{x}(5)$-$\mathbf{x}(8)$, and this procedure continues until $\mathbf{x}(21)$-$\mathbf{x}(24)$ are built by diagonal shift of the transmission vectors of the previous round. In total, $6 \times 4 = 24$ transmission intervals are required and each missing part will appear $2 + 3 = 5$ times; resulting

---

**Algorithm 1** $\mathbf{R}_1$ Generation Procedure

1: **procedure** GENERATE $\mathbf{R}_1$
2:     **for all** $j \in [1 : t]$ **do**
3:        **for all** $i \in [1 : K - 2t + j]$ **do**
4:           $\mathbf{R}_1[i, j] \leftarrow t + i$
5:        **for all** $i \in [K - 2t + j + 1 : K - t]$ **do**
6:           $\mathbf{R}_1[i, j] \leftarrow t + i - (K - 2t + j - 1) - t$
7:     **for all** $j \in [t + 1 : t + L]$, $i \in [1 : K - t]$ **do**
8:        $\mathbf{R}_1[i, j] \leftarrow 1$

---

in total subpacketization requirement of $6 \times 5 = 30$.

For a general network setup with parameters $K, t, L$, in a single round we have $K - t$ transmission vectors, and each new round is built by diagonal shift of transmission vectors in the previous round. There exist a total number of $K$ rounds, resulting in $S = K \times (K - t)$ total transmission intervals; and the required subpacketization is $K \times (t + L)$.

### B. Delivery Prime Matrices

In order to provide the delivery algorithm, we first introduce and construct Delivery Prime (DP) matrices. Denoted by $\mathbf{R}_k$ and $\mathbf{C}_k$, $k \in [K]$, they are a group of matrices with dimensions $(K - t) \times (t + L)$. $\mathbf{R}_1$ and $\mathbf{C}_1$ are built using Algorithms 1 and 2 (for all algorithms, $K, L, N, t$ are assumed to be global variables known to the procedures); and for $k > 1$, we use circular increment to build $\mathbf{R}_k$ and $\mathbf{C}_k$ from $\mathbf{R}_{k-1}$ and $\mathbf{C}_{k-1}$, respectively. For an integer $a$, the intended circular increment operation in domain $K$ is defined as

$$
i_c(a, K) = (a \mod K) + 1 \;,
\tag{11}
$$

while for a matrix $\mathbf{A}$ with positive integer elements, $i_c(\mathbf{A}, K)$ results in a matrix in which each element is the circular increment (in domain $K$) of its respective element in $\mathbf{A}$. Now for $k \in [2 : K]$ we define

$$
\mathbf{R}_k = i_c(\mathbf{R}_{k-1}, K) \;; \quad \mathbf{C}_k = i_c(\mathbf{C}_{k-1}, K) \;.
\tag{12}
$$

**Example 2.** *For the network of Example 1, we have*

$$
\mathbf{R}_1 = \begin{bmatrix} 3 & 3 & 1 & 1 & 1 \\ 4 & 4 & 1 & 1 & 1 \\ 5 & 5 & 1 & 1 & 1 \\ 2 & 6 & 1 & 1 & 1 \end{bmatrix}, \; \mathbf{C}_1 = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 1 & 2 & 4 & 5 & 6 \\ 1 & 2 & 5 & 6 & 3 \\ 1 & 2 & 6 & 3 & 4 \end{bmatrix},
$$

$$
\mathbf{R}_2 = \begin{bmatrix} 4 & 4 & 2 & 2 & 2 \\ 5 & 5 & 2 & 2 & 2 \\ 6 & 6 & 2 & 2 & 2 \\ 3 & 1 & 2 & 2 & 2 \end{bmatrix}, \; \mathbf{C}_2 = \begin{bmatrix} 2 & 3 & 4 & 5 & 6 \\ 2 & 3 & 5 & 6 & 1 \\ 2 & 3 & 6 & 1 & 4 \\ 2 & 3 & 1 & 4 & 5 \end{bmatrix}.
\tag{13}
$$

### C. The Delivery Algorithm

Delivery procedure is provided in Algorithm 3; with its auxiliary procedures presented in Algorithms 4 and 5. We have used $\{\mathbf{C}_k\} \equiv \{\mathbf{C}_1, ..., \mathbf{C}_K\}$ and $\{\mathbf{R}_k\} \equiv \{\mathbf{R}_1, ..., \mathbf{R}_K\}$; and MODULATE function returns the modulated version of its input. The procedure is based on DP matrices. For each $k \in [K]$ we use $\mathbf{R}_k$ and $\mathbf{C}_k$ jointly, to create one transmission round. Each round has $K - t$ transmission vectors; resulting in $S = K \times (K - t)$ transmission intervals in total.

As an explanation, for every $k \in [K]$, a transmission vector $\mathbf{x}(s)$ is built for each row in $\mathbf{R}_k$ and $\mathbf{C}_k$. More precisely,

**Algorithm 2** $C_1$ Generation Procedure

1: **procedure** GENERATE $C_1$
2:     **for all** $j \in [1:t]$, $i \in [1:K-t]$ **do**
3:         $C_1[i,j] \leftarrow j$
4:     **for all** $j \in [t+1:t+L]$, $i \in [1:K-t]$ **do**
5:         **if** $j+i-1 \leq K$ **then**
6:             $C_1[i,j] \leftarrow j+i-1$
7:         **else**
8:             $C_1[i,j] \leftarrow j+i-1-(K-t)$

---

**Algorithm 3** Delivery Procedure

1: **procedure** DELIVERY($\{R_k\}, \{C_k\}$)
2:     INITIALIZE
3:     **for all** $k \in [K]$ **do**
4:         **for all** $i \in [K-t]$ **do**
5:             $s \leftarrow s+1$
6:             $\mathbf{x}(s) \leftarrow 0$
7:             **for all** $j \in [t+L]$ **do**
8:                 $usr \leftarrow C_k[i,j]$
9:                 $prt \leftarrow R_k[i,j]$
10:                $ind \leftarrow q\big(W(usr), prt\big)$
11:                $\mathsf{X} \leftarrow$ MODULATE$(W_{prt}^{ind}(usr))$
12:                $\mathcal{R} \leftarrow$ GENERATE $\mathcal{R}(k,i,usr,prt,\{C_k\})$
13:                $\mathbf{x}(s) \leftarrow \mathbf{x}(s) + \mathbf{v}_{\mathcal{R}}\mathsf{X}$
14:                $q\big(W(usr), prt\big) \leftarrow q\big(W(usr), prt\big) + 1$
15:         Transmit $\mathbf{x}(s)$

---

the entries in each row of $C_k$ specify the user indices to be targeted during the transmission interval, i.e. $\mathcal{T}(s)$; while the entries in the respective row of $R_k$ clarify the file parts to be selected for users in $\mathcal{T}(s)$. Finally, for $n \in [N]$, $k \in [K]$, the superscript $q(n,k)$ is another index which exactly specifies which data portion should be sent during each transmission.

## VI. VALIDITY AND PERFORMANCE

**Lemma 1.** *During each transmission interval, exactly $t+L$ users receive part of their requested data, interference free.*

*Proof.* According to graphical and algorithmic representations of Section V, the delivery phase consists of $K$ rounds; and each round includes $K-t$ transmission vectors. However, the transmission vectors at each round are diagonal shifts of the ones at the previous round, and cache placement structure is also immune to the diagonal shift operation. As a result, it is enough to show that all transmission vectors at a single round serve $t+L$ users interference free.

Without loss of generality, let us consider the first round. Using $\mathcal{T}_r^j$ to denote the set of users targeted at transmission $j$ of round $r$, we have $\mathcal{T}_1^1 = [1:t+L]$. We distinguish two disjoint subsets of $\mathcal{T}_1^1$; as $\mathcal{V}_1^1 = [1:t]$ and $\mathcal{H}_1^1 = [t+1:t+L]$. For any $v \in \mathcal{V}_1^1$, we transmit $W_t(v)$ (superscript $q$ is ignored for simplicity), which is available at the cache memory of $t$ users in $\mathcal{H}_1^1$. So, for interference free data delivery, $W_t(v)$ should be nulled out at the other $L-t$ users of $\mathcal{H}_1^1$; as well as all $t-1$ users of $\mathcal{V}_1^1 \setminus \{v\}$. On the other hand, for any $h \in \mathcal{H}_1^1$,

**Algorithm 4** Initialization Procedure

1: **procedure** INITIALIZE
2:     $s \leftarrow 0$
3:     **for all** $n \in [N]$, $p \in [K]$ **do**
4:         $q(n,p) \leftarrow 1$

---

**Algorithm 5** $\mathcal{R}$ Generation Function

1: **function** GENERATE $\mathcal{R}(k,row,usr,prt,\{C_k\})$
2:     $\mathcal{R} \leftarrow \{usr\}$
3:     **for all** $r \in [t+L]$ **do**
4:         $node \leftarrow C_k[row,r]$
5:         **if** $V[prt,node] = 1$ **then**
6:             $\mathcal{R} \leftarrow \mathcal{R} \cup \{node\}$
7:     **return** $\mathcal{R}$

---

$W_1(h)$ is available at the cache memory of all users in $\mathcal{V}_1^1$; and should be nulled out at all $L-1$ users of $\mathcal{H}_1^1 \setminus \{h\}$. So every data part in the transmission vector should be nulled out at $L-1$ users in total, which is indeed possible with $L$ transmitting antennas. This means all users in $\mathcal{T}_1^1$ can get part of their requested data, interference free.

Subsequent transmission vectors in round one are designed by circular shift of users in $\mathcal{V}_1^1$ and $\mathcal{H}_1^1$, in vertical and horizontal directions respectively. However, for any $j \in [K-t]$, the file part intended for each user in $\mathcal{H}_1^j$ is always available in the cache memory of all users in $\mathcal{V}_1^j$; and should be nulled out at $L-1$ other users of $\mathcal{H}_1^j$. Moreover, the diagonal structure of the cache placement causes the file part requested by each user in $\mathcal{V}_1^j$ to be available in the cache memories of $t$ users in $\mathcal{H}_1^j$; which means it should be also nulled out at $L-1$ total users. So every transmission vector at round one delivers data to $t+L$ users interference free, and the proof is complete. $\quad\square$

**Lemma 2.** *For each $W \in \mathcal{F}$, $W_p$ should be split into $t+L$ smaller parts, for the algorithm to work properly.*

*Proof.* Using the same notation as the proof of Lemma 1, in round $r$, transmission vectors are built by circular shift of $\mathcal{V}_r^1$ and $\mathcal{H}_r^1$, in vertical and horizontal directions respectively. As a consequence, during round $r$:
- for each user $k \in \bigcup \mathcal{H}_r^j$, $W_r(k)$ appears $L$ times in the transmission vectors;
- for each user $k \in \mathcal{V}_r^1$ and $p \in [K]$ such that $V[p,k]=0$, $W_p(k)$ appears once in the transmission vectors.

Moreover, for any $k,p \in [K]$ with $V[p,k]=0$, there exists exactly one $r$ value for which $W_p(k)$ appears in $\bigcup \mathcal{H}_r^j$; but there exist $t$ different $r$ values for which $W_p(k)$ appears in $\mathcal{V}_r^j$ (for some $j$). So $W_p(k)$ appears $L \times 1 + 1 \times t = L+t$ times during all transmissions; which means each $W_p$ should be split in $t+L$ smaller parts for the algorithm to work properly. $\quad\square$

**Corollary 1.** *The provided coded caching scheme achieves the maximum possible DoF of $t+L$ during each transmission and requires linear subpacketization of $K \times (t+L)$.*

Linear growth in subpacketization enables coded caching to be practically implemented in large networks. For example, if
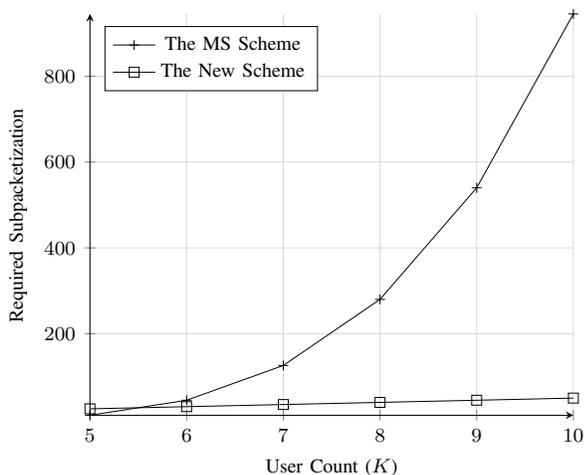
Fig. 3: Subpacketization Comparison - $t = 2$, $L = 3$

$K = 20$, $L = 4$, $t = 2$, the multi-server (MS) scheme of [11] requires subpacketization of 129,200; while our algorithm reduces it to 120 (1,000 fold decrease). Keeping $L$ and $t$ fixed and increasing $K$ to 50, the reduction becomes 66,000 fold.

For very small networks however, there exist specific cases where the MS scheme requires smaller subpacketization. For example, for $t = 2$, $L = 3$, $K \in [5 : 10]$, subpacketization of both schemes is plotted in Figure 3. Clearly, the MS scheme outperforms the new scheme at $K = 5$. However, our scheme has smaller subpacketization as $K$ is increased, and the gap between the two schemes grows exponentially.

Compared with the scheme of [8], as it has DoF loss if either $\frac{K}{L}$ or $\frac{t}{L}$ is non-integer, it is only comparable with our scheme if $t = L$ and at the same time, $\frac{K}{L}$ is an integer. For this special case, it requires subpacketization $\frac{K}{L}$, outperforming our scheme by a factor of $L^3$. However, as mentioned, this is only valid for a very specific class of network parameters and otherwise, the two schemes cannot be directly compared.

## VII. CONCLUSION AND FUTURE WORK

We proposed a coded caching scheme with linear subpacketization; which is applicable for any set of network parameters as long as the the multi-antenna gain is larger than or equal to the global caching gain. Moreover, it achieves the full additive gain of coded caching and multi-antenna communication, in every transmission interval.

The DoF analysis provided in this paper is applicable only to the high-SNR regime, however. It remains an open problem to analyze the system behavior at finite-SNR, where beamformer design complexity is another issue to be considered alongside the large subpacketization requirement.

Another question is the possibility of constructing the DP matrices (and designing the delivery algorithm respectively), using other values of $P$, i.e. $P \neq K$, for the placement matrix. In fact, increasing $P$ might enable better multicasting opportunity, resulting in increased efficiency index and hence better finite-SNR rate, as outlined in [21].

Finally, extending the provided scheme to be applicable to a larger selection of network parameters, i.e. the region $t > L$;

and improving its subpacketization requirement for $t$ values scaling with $K$, is part of our ongoing research.

## REFERENCES

[1] V. N. I. Cisco, "Cisco visual networking index: Forecast and trends, 2017–2022," *White Paper*, vol. 1, 2018.

[2] 6Genesis, "Key Drivers and Research Challenges for 6G Ubiquitous Wireless Intelligence," *White Paper*, vol. 1, 2019.

[3] M. Katz, M. Matinmikko-Blue, and M. Latva-Aho, "6Genesis Flagship Program: Building the Bridges Towards 6G-Enabled Wireless Smart Society and Ecosystem," in *2018 IEEE 10th Latin-American Conference on Communications (LATINCOM)*. IEEE, 2018, pp. 1–9.

[4] F. Boccardi, R. W. Heath, A. Lozano, T. L. Marzetta, and P. Popovski, "Five disruptive technology directions for 5G," *IEEE communications magazine*, vol. 52, no. 2, pp. 74–80, 2014.

[5] A. Osseiran, F. Boccardi, V. Braun, K. Kusume, P. Marsch, M. Maternia, O. Queseth, M. Schellmann, H. Schotten, H. Taoka, and others, "Scenarios for 5G mobile and wireless communications: the vision of the METIS project," *IEEE communications magazine*, vol. 52, no. 5, pp. 26–35, 2014.

[6] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Transactions on Information Theory*, vol. 60, no. 5, pp. 2856–2867, 2014.

[7] A. Leventhal, "Flash storage memory," *Communications of the ACM*, vol. 51, no. 7, pp. 47–51, 2008.

[8] E. Lampiris and P. Elia, "Adding transmitters dramatically boosts coded-caching gains for finite file sizes," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 6, pp. 1176–1188, 2018.

[9] M. A. Maddah-Ali and U. Niesen, "Decentralized coded caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Transactions on Networking (TON)*, vol. 23, no. 4, pp. 1029–1040, 2015.

[10] N. Karamchandani, U. Niesen, M. A. Maddah-Ali, and S. N. Diggavi, "Hierarchical coded caching," *IEEE Transactions on Information Theory*, vol. 62, no. 6, pp. 3212–3229, 2016.

[11] S. P. Shariatpanahi, S. A. Motahari, and B. H. Khalaj, "Multi-server coded caching," *IEEE Transactions on Information Theory*, vol. 62, no. 12, pp. 7253–7271, 2016.

[12] S. P. Shariatpanahi, G. Caire, and B. H. Khalaj, "Multi-antenna coded caching," in *2017 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2017, pp. 2113–2117.

[13] ——, "Physical-layer schemes for wireless coded caching," *IEEE Transactions on Information Theory*, vol. 65, no. 5, pp. 2792–2807, 2018.

[14] A. Tölli, S. P. Shariatpanahi, J. Kaleva, and B. Khalaj, "Multicast beamformer design for coded caching," in *2018 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2018, pp. 1914–1918.

[15] ——, "Multi-antenna interference management for coded caching," *arXiv preprint arXiv:1711.03364*, 2017.

[16] K. Shanmugam, M. Ji, A. M. Tulino, J. Llorca, and A. G. Dimakis, "Finite-length analysis of caching-aided coded multicasting," *IEEE Transactions on Information Theory*, vol. 62, no. 10, pp. 5524–5537, 2016.

[17] Q. Yan, M. Cheng, X. Tang, and Q. Chen, "On the placement delivery array design for centralized coded caching scheme," *IEEE Transactions on Information Theory*, vol. 63, no. 9, pp. 5821–5833, 2017.

[18] Q. Yan, X. Tang, Q. Chen, and M. Cheng, "Placement delivery array design through strong edge coloring of bipartite graphs," *IEEE Communications Letters*, vol. 22, no. 2, pp. 236–239, 2017.

[19] C. Shangguan, Y. Zhang, and G. Ge, "Centralized coded caching schemes: A hypergraph theoretical approach," *IEEE Transactions on Information Theory*, vol. 64, no. 8, pp. 5755–5766, 2018.

[20] K. Shanmugam, A. M. Tulino, and A. G. Dimakis, "Coded caching with linear subpacketization is possible using Ruzsa-Szeméredi graphs," in *2017 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2017, pp. 1237–1241.

[21] M. Salehi, A. Tölli, S. P. Shariatpanahi, and J. Kaleva, "Subpacketization-Rate Trade-off in Multi-Antenna Coded Caching," in *2019 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2019, pp. 1–6.

[22] E. Lampiris and P. Elia, "Bridging two extremes: Multi-antenna Coded Caching with Reduced Subpacketization and CSIT," in *2019 IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, 2019, pp. 1–5.