

Multi-Antenna Coded Content Delivery with Caching: A Low-Complexity Solution

Junlin Zhao, *Member, IEEE*, Mohammad Mohammadi Amiri, *Member, IEEE*, and Deniz Gündüz, *Senior Member, IEEE*

Abstract—We study downlink beamforming in a single-cell network with a multi-antenna base station serving cache-enabled users. Assuming a library of files with a common rate, we formulate the minimum transmit power with proactive caching and coded delivery as a non-convex optimization problem. While this multiple multicast problem can be efficiently solved by successive convex approximation (SCA), the complexity of the problem grows exponentially with the number of subfiles delivered to each user in each time slot, which itself grows exponentially with the number of users. We introduce a low-complexity alternative through time-sharing that limits the number of subfiles received by a user in each time slot. We then consider the joint design of beamforming and content delivery with sparsity constraints to limit the number of subfiles received by a user in each time slot. Numerical simulations show that the low-complexity scheme has only a small performance gap to that obtained by solving the joint problem with sparsity constraints, and outperforms state-of-the-art results at all signal-to-noise ratio (SNR) and rate values with a sufficient number of transmit antennas. A lower bound on the achievable degrees-of-freedom (DoF) of the low-complexity scheme is derived to characterize its performance in the high SNR regime.

I. INTRODUCTION

The seminal work of Maddah-Ali and Niesen showed that by exploiting caches at the users in order to create and exploit multicasting opportunities we can reduce the delivery time, or equivalently increase the throughput in wireless networks [2]. With coded caching, uncoded contents can be proactively pushed at user devices without knowing users' demands, and a server can serve multiple users simultaneously by broadcasting specially designed coded combinations of the

remaining parts of all the users' requests, to guarantee that all the users can recover their desired contents. This feature is particularly favorable in wireless medium due to its broadcast nature. However, [2] ignored the physical characteristics of the channel, and simply assumed error-free communication, and focused on minimizing the number of bits delivered error-free over this multicast channel.

Over the last years many follow-up works have studied coded delivery over noisy broadcast channels. When users may have different channel capacities, the user with the worst channel condition becomes the bottleneck limiting the performance of multicasting. The global caching gain promised in [2] is hence not straightforward in practice. Coded caching in erasure broadcast channels is studied in [3] and [4] by allocating cache memories at weak receivers to overcome this bottleneck. A simple binary Gaussian broadcast channel is considered in [5], and an interference enhancement scheme is used to overcome the limitation of weak users. A cache-aided multicasting strategy over a Gaussian broadcast channel is presented in [6], with superposition coding and power allocation. The authors in [7] consider fading channels, and show that a linear increase in the sum delivery rate with the number of users can be achieved with user selection.

Another important line of research has focused on evaluating the performance of coded caching and delivery in the presence of multiple transmit antennas at the server. Multicast beamforming, where the multiple-antenna base station (BS) multicasts distinct data streams to multiple user groups, is an efficient physical layer technique [8]–[10]. In [11], the authors extend the results in [7] to multi-input single-output (MISO) fading channels, where the same linear increase in content delivery rate with respect to (w.r.t.) the number of users is achieved without channel state information at the transmitter (CSIT), and an improvement is obtained with spatial multiplexing when CSIT is available. In [12], coded delivery is employed along with zero-forcing to simultaneously exploit spatial multiplexing and caching gains. With multiple antennas at the BS, coded messages can be nulled at unintended user groups, which increases the number of users simultaneously served as compared to the single antenna setting. Particularly, this approach was found to achieve the near-optimal degrees-of-freedom (DoF) in [13]. In addition to the gain in content delivery rate, employing multiple transmit antennas also allows reducing the subpacketization level required in coded caching [14], [15].

By treating the transmission of coded subfiles as a coordinated beamforming problem, improved spectral efficiency

Part of this work was presented at the IEEE International Workshop on Signal Processing Advances in Wireless Communications (SPAWC), Cannes, France, Jul. 2019 [1].

This work was partially supported by the European Research Council (ERC) through project BEACON (No. 677854), and by the European Union's Horizon 2020 Research and Innovation Programme through project SCAVENGE (No. 675891).

J. Zhao was with the Information Processing and Communications Lab, Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, UK. He is now with the School of Science and Engineering, the Chinese University of Hong Kong, Shenzhen, Guangdong 518172, P.R. China, and also with the University of Science and Technology of China, Anhui 230026, P.R. China (e-mail: j.zhao15@imperial.ac.uk)

M. Amiri was with the Information Processing and Communications Lab, Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, UK. He is now with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544, USA (e-mail: mamiri@princeton.edu)

D. Gündüz is with the Information Processing and Communications Lab, Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, UK (e-mail: d.gunduz@imperial.ac.uk)

is achieved in [16] by optimizing the beamforming vectors, which is also shown to achieve the same DoF as in [12] in special cases. Memory-sharing is proposed in [17] to apply the content placement scheme of [2] for a fraction of the library, which exploits both the spatial multiplexing gain and the global caching gain by sending a common message together with user-dependent messages. The impact of imperfect CSIT on achievable DoF is considered for MISO broadcast channels in [18]. Similarly to [17], [19] adopts memory-sharing, and proposes a joint unicast and multicast beamforming approach.

In this paper, motivated by the results in [12] and [16], we consider a cache-aided MISO broadcast channel. Firstly, a general framework for cache-aided downlink beamforming is formulated, focusing on the minimum required transmit power for delivering the contents at a prescribed common rate. The resultant nonconvex optimization problem is tackled by successive convex approximation (SCA), which is guaranteed to converge to a stationary solution of the original nonconvex problem. As noted in [16], the beamforming design involves solving an optimization problem with exponentially increasing number of constraints with the number of coded messages each user decodes in each time slot. To limit the complexity, we propose a novel content delivery scheme, in which the coded subfiles, each targeted at a different subset of receivers, are delivered over multiple orthogonal time slots, while the number of coded messages each user decodes in each time slot can be flexibly adjusted. Unlike the scheme in [16], the scheme we propose does not limit the number of users served in each time slot, but directly limits the number of messages each user decodes, and hence, the complexity of the decoder. We propose a greedy algorithm that decides the multicast messages to be delivered at each time slot, and the number of time slots. A lower bound on the DoF achieved by the greedy scheme is also provided. We then consider a more general design of the beamforming vectors together with the content delivery scheme with a constraint on the maximum number of messages each user can decode in any time slot. We formulate this joint optimization as a power minimization problem with sparsity constraints, and solve it via SCA to obtain a stationary solution. Our numerical results show that the proposed greedy scheme has only a small performance gap to that of the optimization-based delivery scheme, and provide significant gains over the one proposed in [16] in terms of transmit power, particularly in the high rate/high signal-to-noise ratio (SNR) regime.

The remainder of the paper is organized as follows. Section II introduces the system model. In Section III, we present an achievable coded delivery scheme, and formulate the power minimization problem for a multi-antenna server. We introduce a low-complexity content delivery scheme in Section IV, and present its DoF analysis in Section V. In Section VI, we consider the joint design of beamforming and coded content delivery. Finally, we compare the proposed schemes with the state-of-the-art through numerical simulations in Section VII, and conclude the paper in Section VIII.

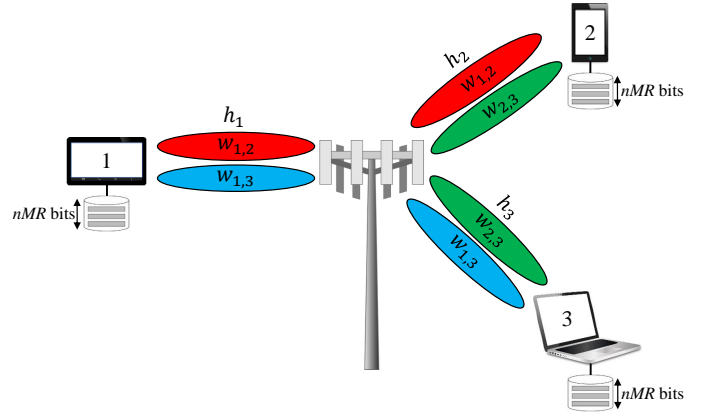


Figure 1: Illustration of a cache-aided MISO channel with $K = 3$ users. Multi-antenna BS employs multicast beamforming to deliver the missing parts of users' requests.

II. SYSTEM MODEL

We consider downlink transmission within a single cell, where a BS equipped with N_T antennas serves K single-antenna cache-equipped users, as illustrated in Fig. 1. We consider a library of N files, denoted by \mathcal{V} , $(V_1; \dots; V_N)$, each distributed uniformly over the set 2^{nR} , available at the BS, where R and n represent the rate of each file and the blocklength, respectively. Each user is equipped with a local cache that can store up to M files, and the corresponding *caching factor*, t , is defined as the ratio of the total cache capacity across all the receivers to the library size, $t = MK/N$.

Contents are placed at users' caches during off-peak periods without any prior information on the user requests or the channel state information (CSI) to be experienced during the delivery phase. Caching function for user k is denoted by $\binom{(n)}{k} : 2^{nR} \times N \rightarrow 2^{nMR}$, which maps the library to the cache contents Z_k at user k , i.e., $Z_k = \binom{(n)}{k}(\mathcal{V})$, $k \in [K]$. Once the users reveal their demands \mathbf{d} , $(d_1; \dots; d_K)$, where $d_k \in [N]$; $8k \in [K]$, signal $\mathbf{X} \in \mathbb{C}^{N_T \times n}$ is transmitted, where $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]$, and $\mathbf{x}_i \in \mathbb{C}^{N_T \times 1}$ is the channel input vector at time i ; $i = 1; \dots; n$. An average power constraint P is imposed on each channel input \mathbf{X} . User k receives

$$\mathbf{y}_k^T = \mathbf{h}_k^H \mathbf{X} + \mathbf{n}_k^T; \quad (1)$$

where $\mathbf{h}_k \in \mathbb{C}^{N_T \times 1}$ is the channel vector from the BS to the k -th user, and $\mathbf{n}_k \in \mathbb{C}^{n \times 1}$ is the additive white Gaussian noise at user k with each entry independent and identically (i.i.d.) distributed according to $\mathcal{CN}(0; \frac{1}{2})$, $k \in [K]$. We assume that the CSI is perfectly known to the BS and the receivers in the delivery phase. Hence, the encoding function at the BS, $\binom{(n)}{k} : 2^{nR} \times N \times [N]^K \rightarrow \mathbb{C}^{N_T \times n}$, maps the library, the demand vector, and the CSI to the channel input vector. We note here that, while the channel encoding function $\binom{(n)}{k}$ depends on the demand vector and the CSI, caching functions $\binom{(n)}{k}$ depend only on the library. After receiving \mathbf{y}_k , user k reconstructs \hat{V}_{d_k} using its local cache

¹For any positive real X , we define $\lfloor X \rfloor$ as the set of positive integers less than or equal to X .

content Z_k , channel vector \mathbf{h}_k , and demand vector \mathbf{d} through function $\hat{V}_k^{(n)} : \mathbb{C}^n \times \mathbb{R}^{2nMR} \times \mathbb{C}^{N_T-1} \times [N]^K \rightarrow \mathbb{R}^{2nR}$, i.e., $\hat{V}_k = \hat{V}_k^{(n)}(\mathbf{y}_k; Z_k; \mathbf{h}_k; \mathbf{d})$, $k \in [K]$. The probability of error is defined as $P_e = \max_{\mathbf{d}} \max_{k \in [K]} \Pr\{V_{d_k} \notin \hat{V}_k\}$. An $(R; M; P)$ tuple is *achievable* if there exist a sequence of caching functions $\hat{V}_1^{(n)}; \dots; \hat{V}_K^{(n)}$, encoding function $\mathbf{e}^{(n)}$, and decoding functions $\hat{V}_1^{(n)}; \dots; \hat{V}_K^{(n)}$, such that $P_e \rightarrow 0$ as $n \rightarrow \infty$. For file rate R and cache size M , our goal is to characterize

$$P(R; M) = \inf \{P : (R; M; P) \text{ is achievable}\} \quad (2)$$

which characterizes the minimum required transmit power that guarantees the reliable delivery of any demand vector.

III. AN ACHIEVABLE DELIVERY SCHEME

In this section, we present a multi-antenna transmission scheme with coded caching, where the cache placement and coded content generation follows [2], while beamforming is employed at the BS to multicast coded subfiles to the receivers.

A. Placement and Delivery Schemes

For a caching factor $t \in \{1, \dots, K-1\}$, we represent t -element subsets of $[K]$ by $G_1^t; \dots; G_{\binom{K}{t}}^t$. File V_i , $i \in [N]$, is divided equally into $\binom{K}{t}$ disjoint subfiles $V_{i;G_1^t}; \dots; V_{i;G_{\binom{K}{t}}^t}$, each consisting of $n \frac{R}{\binom{K}{t}}$ bits. User k , $k \in [K]$, caches subfile $V_{i;G_j^t}$, if $k \in G_j^t$, $j \in [\binom{K}{t}]$. The cache content of user k is then given by $\bigcup_{i \in [N]} \bigcup_{j \in [\binom{K}{t}]: k \in G_j^t} V_{i;G_j^t}$.

During the *delivery phase*, for any demand combination \mathbf{d} , we aim to deliver the coded message

$$S_{G_j^{t+1}} = \bigcup_{k \in G_j^{t+1}} V_{d_k; G_j^{t+1}} \quad (3)$$

to all the users in set G_j^{t+1} , for $j \in [\binom{K}{t+1}]$. Observe that, after receiving $S_{G_j^{t+1}}$, each user $k \in G_j^{t+1}$ can recover subfile $V_{d_k; G_j^{t+1}}$ having access to $V_{d_i; G_j^{t+1}}$, $\forall i \in G_j^{t+1}$.

We define $S = \{G_1^{t+1}; \dots; G_{\binom{K}{t+1}}^{t+1}\}$ as the set of all the multicast messages, with each message $T \in S$ represented by the set of users it is targeting, and let $S_k \subseteq S$ denote the subset of messages targeting user k . We have $|S_k| = \binom{K}{t+1}$ and $|S_k \cap S_j| = \binom{K}{t+1}$.

The following examples will be used to better explain the proposed scheme:

Example 1: Let $N = 5$; $K = 5$; $M = 1$. We have $t = \frac{MK}{N} = 1$. Each file is split into $\binom{K}{t} = 5$ disjoint subfiles of the same size, where we represent file i , $i \in [N]$, as

$$V_i = [V_{i;f1g}; V_{i;f2g}; V_{i;f3g}; V_{i;f4g}; V_{i;f5g}] \quad (4)$$

The cache content of user k is $Z_k = [V_{i;fk g}, k \in [K]]$, which satisfies the cache capacity constraint. All user demands

can be fulfilled by delivering the following $\binom{K}{t+1} = 10$ subfiles:

$$\begin{aligned} S_{f1;2g} &= V_{d1;f2g} & V_{d2;f1g}; & S_{f1;3g} &= V_{d1;f3g} & V_{d3;f1g}; \\ S_{f1;4g} &= V_{d1;f4g} & V_{d4;f1g}; & S_{f1;5g} &= V_{d1;f5g} & V_{d5;f1g}; \\ S_{f2;3g} &= V_{d2;f3g} & V_{d3;f2g}; & S_{f2;4g} &= V_{d2;f4g} & V_{d4;f2g}; \\ S_{f2;5g} &= V_{d2;f5g} & V_{d5;f2g}; & S_{f3;4g} &= V_{d3;f4g} & V_{d4;f3g}; \\ S_{f3;5g} &= V_{d3;f5g} & V_{d5;f3g}; & S_{f4;5g} &= V_{d4;f5g} & V_{d5;f4g}; \end{aligned}$$

Example 2: Let $N = 4$; $K = 4$; $M = 1$. We have $t = \frac{MK}{N} = 1$. Each file is split into $\binom{K}{t} = 4$ disjoint subfiles of the same size. All user demands can be fulfilled by delivering the following $\binom{K}{t+1} = 6$ subfiles:

$$S_{f1;2g}; S_{f1;3g}; S_{f1;4g}; S_{f2;3g}; S_{f2;4g}; S_{f3;4g} \quad (5)$$

Note that the message S_T is intended for users in set T , but interferes with users in set $[K] \setminus T$. Moreover, for any demand combination \mathbf{d} , all the users are required to decode the same number of messages, which is $\binom{K}{t+1}$.

B. Multi-Antenna Transmission Scheme

The delivery of the coded messages in set S to their respective receivers is a multi-antenna multi-message multicasting problem. Before introducing our low-complexity scheme in the next section, we present here a general transmission strategy based on message-splitting and time-division transmission. The messages in S can be transmitted over B orthogonal time slots, the i -th of which is of blocklength n_i ; $i \in [B]$, where $\sum_{i=1}^B n_i = n$. The transmitted signal $\mathbf{X}(i) = [\mathbf{x}^P_{j=1}^{n_j+1}; \dots; \mathbf{x}^P_{j=1}^{n_j}]$ at time slot $i \in [B]$ is given by

$$\mathbf{X}(i) = \sum_{T \subseteq S} \mathbf{w}_T(i) S_T^T(i); \quad (6)$$

where $S_T^T(i) \in \mathbb{C}^{n_i-1}$ is a unit-power complex Gaussian signal of blocklength n_i , modulated from the corresponding message S_T in (3), intended for the users in set T , transmitted in time slot i , encoded by the beamforming vector $\mathbf{w}_T(i) \in \mathbb{C}^{N_T-1}$.

The received signal at user k in time slot i is $\mathbf{y}_k^T(i) = \mathbf{h}_k^H \sum_{T \subseteq S} \mathbf{w}_T(i) S_T^T(i) + \mathbf{h}_k^H \sum_{T \subseteq S_k^c} \mathbf{w}_T(i) S_T^T(i) + \mathbf{n}_k^T(i)$; $\underbrace{\sum_{T \subseteq S_k} \{S_T^T(i)\}}_{\text{desired messages}} \quad \underbrace{\sum_{T \subseteq S_k^c} \{S_T^T(i)\}}_{\text{interference}}$ (7)

where S_k^c is the complement of set S_k in S . Let \mathcal{S}_k denote the collection of all non-empty subsets of S_k , with each element of \mathcal{S}_k denoted by $J_{S_k}; j \in [2^{\binom{K}{t+1}} - 1]$. We denote $S(i) \subseteq S$ as the subset of messages transmitted in time slot i , i.e., $T \in S(i)$ if $\mathbf{w}_T(i) \neq \mathbf{0}$.

Note that each user may receive more than one message in each transmission slot. From the capacity region of the associated Gaussian multiple access channel, following conditions must be satisfied for successful decoding of all the intended messages at user k , $k \in [K]$, at time slot i :

$$\sum_{T \in \mathcal{S}_k} R^T(i) \leq \frac{n_i}{n} \log_2 \left(1 + \sum_{T \in \mathcal{S}_k^c} \frac{P_T(i)}{P_k(i)} \right); \quad \forall J_{S_k} \in \mathcal{S}_k; \quad (8)$$

where $R^T(i)$ is the rate of message $s_T(i)$, and $\gamma_k^T(i)$ is the received signal-to-interference-plus-noise ratio (SINR) of message $s_T(i)$ at user k at time slot i , given by

$$\gamma_k^T(i) = \frac{p_j \mathbf{h}_k^H \mathbf{w}_T(i) j^2}{\sum_{l \in S(i)} p_l \mathbf{h}_k^H \mathbf{w}_l(i) j^2 + \frac{\sigma^2}{k}}; \quad (9)$$

for any $T \in \mathcal{K}$, or equivalently, any $T \in \mathcal{S}_k$. The rate of message T is the sum of the rate of submessages $s_T(i)$, and must satisfy

$$\sum_{i=1}^B R^T(i) \leq \frac{R}{K}; \quad \forall T; \quad (10)$$

Note that this scheme is quite flexible; each multicast message can be split into B messages and transmitted over B time slots. It can be specialized to different content delivery schemes by specifying the subset of transmitted subfiles in each time slot and the blocklength of each time slot, i.e., $fS(i)g_{i=1}^B$ and $f n_i g$. Let

$$v_T(i) = \begin{cases} 1 & \text{if } T \in S(i) \\ 0 & \text{if } T \notin S(i) \end{cases} \quad (11)$$

be the indicator function specifying whether message T is transmitted at time slot i or not. Note that $\sum_{i=1}^B v_T(i) = 1$ is required to fulfill users' demands, where $v_T = [v_T(1) \dots v_T(B)]$. It is readily seen that $v_T(i)$ can be inferred by the corresponding beamforming vector $\mathbf{w}_T(i)$, or equivalently, by the message rate $R^T(i)$.

C. Transmit Power Minimization

For any given delivery scheme specified by $v_T(i)$ and n_i , $\forall i \in [B]; \forall T \in \mathcal{S}$, the associated minimum required transmit power problem is obtained as follows:

$$P_{\min} = \min_{\mathbf{w}_T(i); f R^T(i) g; \forall T \in \mathcal{S}; i=1}^B \sum_{i=1}^B \sum_{k \in \mathcal{S}_k} \frac{n_i}{n} k \mathbf{w}_T(i)^2 \quad (12a)$$

$$\text{s.t.} \quad \sum_{T \in \mathcal{S}_k} R^T(i) \leq \frac{n_i}{n} \log_2 \left(1 + \sum_{T \in \mathcal{S}_k} \gamma_k^T(i) \right); \quad (12b)$$

$$\sum_{i=1}^B R^T(i) \leq \frac{R}{K}; \quad \forall T; \quad (12c)$$

$$R^T(i) = 0; \quad \forall i; \forall T; \quad (12d)$$

$$R^T(i) = 0; \quad \text{if } v_T(i) = 0; \quad \forall i; \quad (12e)$$

where $\gamma_k^T(i)$ is defined in (9). Here, constraints in (12b) guarantee that the rates of the messages targeting each user in each time slot are within the capacity region, constraints in (12c) ensure that sufficient information is delivered for each coded subfile over B time slots, while (12d)-(12e) represent the specific content delivery scheme.

Note that the problem in (12) is a generalization of various well-known NP-hard problems depending on the specific content delivery scheme. For $B = jS(j)$ with $jS(j) = 1; \forall i$, the problem boils down to a series of standard multicast beamforming problems, where a common message is broadcast to a different subset of $t+1$ users in each time slot [8].

When $jS(j) > 1$, $T \in \mathcal{S}^0$; if $T \notin \mathcal{S}^0$, $\forall i$, and $S(i) = S(j)$; if $i \notin j$, we need to solve the conventional multigroup multicast beamforming problem at each time slot [9]. It can be seen from (12) that the content delivery scheme specified by $v_T(i)$ and n_i affects the minimum required power. A straightforward solution (as done in [17], [20]) would transmit a single coded message in each time slot. However, this does not exploit the spatial multiplexing gain provided by multiple antennas, and results in a poor DoF performance in the high SNR regime. Another approach studied in [21] is to deliver the coded messages targeting non-overlapping user groups in parallel. Obviously, the content delivery scheme is an important factor on the system performance and needs to be carefully designed.

We remark here that, even when the delivery scheme is specified, the problem in (12) is computationally intractable due to the non-convex constraints in (12b). However, we show in the Appendix that SCA methods [22] can be employed to obtain a stationary point of the problem, which serves as an upper bound on the optimal solution. Starting with a feasible initial point, the SCA algorithm solves a sequence of subproblems in an iterative manner, where the subproblem in the t -th iteration is derived by convexifying the original problem at the solution point of the $(t-1)$ -th subproblem. More detailed discussions on the SCA algorithm are provided in the Appendix.

IV. A LOW-COMPLEXITY DESIGN

In this section, we propose a low-complexity content delivery scheme with the flexibility to adjust the number of coded messages intended for each user at each time slot. Observing that if a set $S(i) = \{T | v_T(i) = 1\}$ of messages are transmitted in time slot i , $c_k(i) = |S(i) \cap \mathcal{S}_k|$ messages are transmitted to user k , which results in $2^{c_k(i)} - 1$ constraints only for user k in time slot i in problem (12). Computational complexity of problem (12) increases drastically with the number of constraints, rendering the numerical optimization problem practically infeasible. More importantly, a multi-user detection scheme needs to be employed at the users, whose complexity also increases with $c_k(i)$.

A low complexity scheme is proposed in [16] by limiting the number of users to be served in each time slot, thereby indirectly reducing the number of coded messages to be decoded by each user. Specifically, an integer parameter $t \in [\min\{N_T, K\} - t_0]$ is leveraged in [16] to control the number of active users in each time slot, which is set to $t+$, and leads to a content delivery scheme with $B = \frac{K}{t+}$ time slots. In each time slot, a fraction of the desired coded messages for all the active users are transmitted. In addition to $t+$, another integer parameter τ determines the possible set partitions of the user subset in each time slot. When $t+$ is divisible by $\tau+$, the user subset can be partitioned into $\frac{t+}{\tau+}$ non-overlapping subsets, and a fraction of the desired coded messages for each partition can be transmitted simultaneously. It is shown in [16] that the system performance can be improved if multiple groups of messages can be transmitted in parallel, i.e., $\frac{t+}{\tau+} \geq 2$, as compared to the case $\tau = t+$.

Moreover, the number of messages for each user to decode in each time slot is $\binom{t+1}{t}$, which is an exponential function of t . Therefore, by adjusting the value of t , the number of coded messages for each user in each time slot is indirectly adjusted.

Instead of limiting the subsets of users to be served in each time slot, we propose to directly adjust the number of coded messages targeted to each user. We will show that this results in a more efficient delivery scheme than the one in [16]. In Example 1, if we transmit all the messages in one time slot, i.e., $B = 1$, a total of $jSj = \binom{K}{t+1} = 10$ coded subfiles are transmitted simultaneously, with each user decoding $\binom{K}{t} = 4$ messages. Accordingly, in the optimization problem in (12) we will have $K \cdot (2^{jSj} - 1) = 75$ constraints. To alleviate the computational complexity, the low complexity scheme in [16] splits each subfile into 3 minifiles, and the coded messages are grouped to serve a subset of $t + 1 = 3$ users in each of the $B = \binom{K}{t+1} = 10$ time slots. Within each time slot, each user needs to decode 2 messages. Note that the power minimization problem for each time slot can be solved independently; therefore, we would need to solve 10 smaller optimization problems, each with $3 \cdot 3 = 9$ constraints.

In contrast, we propose to serve as many users as needed at each time slot while keeping $c_k(i)$ under a given threshold s for each user k . In our Example 1, we can satisfy all the user requests in only 2 time slots, by setting nonzero rate targets for the messages in

$$S(1) = \{f1; 2g; 3g; f3; 4g; f4; 5g; f1; 5gg\}; \text{ and} \\ S(2) = \{f1; 3g; f2; 4g; f3; 5g; f1; 4g; f2; 5gg\}$$

in time slots 1 and 2, respectively. Note that each user k decodes only $c_k(i) = s = 2$ messages in each time slot, the same as the delivery scheme in [16], requiring the same implementation complexity at each user; however, 5 users are served in each time slot, which results in a significantly smaller number of time slots. Thus, we need to solve only two optimization problems at the BS, each with $5 \cdot 3 = 15$ constraints.

In general, the number of constraints in the optimization problem in (12) increases exponentially with s , which results in exponentially increasing number of constraints in the problem in each SCA iteration. Thus the computational complexity of the delivery scheme can be largely alleviated by choosing a small s value, which also simplifies the multi-user detection algorithm.

The key idea of our proposed low-complexity scheme is to divide set S into disjoint subsets $S(1); \dots; S(B)$, with $c_k(i) \leq s, \forall k; i$, while keeping B as small as possible. Since the total number of subfiles to transmit is fixed, choosing a small value of B , i.e., completing the delivery phase within a small number of time slots, requires multiplexing more messages in each time slot, without increasing the complexity of the receivers. To obtain this low-complexity scheme, the

following optimization problem can be formulated:

$$\min_{f, v_T, g; B} B \quad (13a)$$

$$\text{s.t.} \quad v_T(i) \leq s; \forall i \in [B]; k \in [K]; \quad (13b)$$

$$\sum_{i=1}^{T \cdot 3k} v_T(i) = 1; \forall T; \quad (13c)$$

$$v_T(i) \geq f0; 1g; \forall i \in [B]; \quad (13d)$$

where constraint (13b) imposes that each user decodes no more than s messages in each time slot, while (13c) requires that each message will be transmitted in only one time slot. However, since the problem itself varies with variable B , the problem is not in a tractable form. By introducing $L = B$ as a prescribed parameter that determines the dimension of the problem, and an auxiliary variable $q \geq f0; 1g^L$, problem (13) can be equivalently written as

$$B = \min_{f, v_T, g; q} \mathbf{1}^T \mathbf{q} \quad (14a)$$

$$\text{s.t.} \quad v_T(i) \leq s; \forall i \in [L]; k \in [K]; \quad (14b)$$

$$\sum_{i=1}^{T \cdot 3k} v_T(i) = 1; \forall T; \quad (14c)$$

$$\sum_{i=1}^T v_T(i) \leq \sum_{i=1}^K q_i; \forall i \in [L]; \quad (14d)$$

$$v_T(i) \geq f0; 1g; \forall i \in [L]; \quad (14e)$$

$$q \geq f0; 1g^L; \quad (14f)$$

where $\mathbf{1}$ denotes a column vector of all ones. Since $\frac{K}{t+1}$ is a bound on $\sum_{i=1}^T v_T(i)$, the optimal q_i is 1 if $\sum_{i=1}^T v_T(i)$ is nonzero, and 0 otherwise, in order to minimize $\sum_{i=1}^K q_i$ in the objective. Note that problem (14) can be considered as minimizing the number of time slots employed out of a maximum L available time slots. We can set $L = \frac{K}{t+1}$ which guarantees the existence of a solution; however, choosing a smaller L will reduce the complexity of the problem. The problem in (14) is a 0-1 integer programming problem, which is generally NP-hard [23].

Algorithm 1 Low-complexity greedy delivery scheme

Input: $N; K; M; s; R$
Output: $B, \sum_{i=1}^B fS(i)g, \sum_{i=1}^B fn_i g; \delta T$

- 1: Set $t = \frac{MK}{N}$, $i = 1$, and $E = S$
- 2: **while** $E \neq ?$ **do**
- 3: Set $c(i), [c_1(i) \dots c_K(i)] = \mathbf{0}, S(i) = ?, C = E$
- 4: **while** $c_k(i) + 1 \leq s; \delta k \geq [K]$ and $C \neq ?$ **do**
- 5: $K, \{fkj\} \arg \min_k c(i)g$
- 6: Find $\hat{T} = \arg \max_{T \subseteq C} jK^T Tj$
- 7: $C = C \setminus nf\hat{T}g$
- 8: **if** $c_k(i) + 1 \leq s; \delta k \geq \hat{T}$ **then**
- 9: $c_k(i) = c_k(i) + 1; \delta k \geq \hat{T}$
- 10: $S(i) = S(i) \cup \hat{T}, E = E \setminus nf\hat{T}g$
- 11: **else**
- 12: **break**
- 13: **end if**
- 14: **end while**
- 15: $i = i + 1$
- 16: **end while**
- 17: Set $B = i - 1$
- 18: **for** $i = 1 : B$ **do**
- 19: $n_i = \frac{jS(i)j}{\binom{K+i}{t}} n$
- 20: $R^T(i) = \frac{R}{\binom{K}{t}}; \delta T \geq S(i)$
 0 ; otherwise
- 21: **end for**

In Algorithm 1, we propose a greedy solution that constructs disjoint $S(i)$ sets for any s value. Specifically, $S(i)$'s are generated in a sequential manner: to construct $S(i)$, we initialize $c(i), [c_1(i) \dots c_K(i)] = \mathbf{0}, S(i) = ?$, and the set $E = S \setminus \sum_{j=1}^{i-1} S(j)$ of remaining messages for assignment, we first identify the user(s) that have decoded the least number of messages so far, i.e., user(s) in set $K, \{fkj\} \arg \min_k c(i)g$, and check whether there exists a message $\hat{T} \subseteq E$ such that the condition $c_k(i) + 1 \leq s$ for $\delta k \geq \hat{T}$ holds. If no such \hat{T} can be found, the process of constructing $S(i)$ is completed, and we start constructing $S(i+1)$ in the same manner. The whole procedure is completed when $E = ?$, i.e., all the messages have been assigned to a subset. Note that our proposed greedy scheme covers the case of $B = 1$, where all the messages are sent simultaneously, if $s \geq \frac{K+1}{t}$.

Next we elaborate the proposed greedy content delivery algorithm in Examples 1 and 2.

Example 1 (continued): $N = 5; K = 5; M = 1. t = \frac{MK}{N} = 1$. Suppose $s = 2$. As illustrated in Fig. 2, the algorithm starts by constructing $S(1)$, i.e., identifying the coded messages to be delivered in the first time slot, initialized as $c(1), [c_1(1) \dots c_K(1)] = \mathbf{0}, S(1) = ?$. Firstly, it is obvious that $K = [K]$ since $c_k(1) = 0; \delta k$. Hence, one may choose any of the available messages in $E = S$. Suppose message $S_{f1;2g}$ is chosen. We update $c_1(1) = c_2(1) = 1, E = E \setminus nf1;2g$. The algorithm then identifies the updated $K = [3;4;5]$, according to which one may choose from $S_{f3;4g}; S_{f3;5g}; S_{f4;5g}$ without violating the constraint $c_k(1) + 1 \leq s, \delta k$. Suppose $S_{f3;4g}$ is chosen, and we have $c_1(1) = c_2(1) = c_3(1) =$

$c_4(1) = 1, E = E \setminus nf3;4g$, and $K = [5]$; and accordingly one may choose from $S_{f1;5g}; S_{f2;5g}; S_{f3;5g}; S_{f4;5g}$, while still keeping the constraint $c_k(1) + 1 \leq s, \delta k$. Similarly, messages $S_{f2;3g}$ and $S_{f4;5g}$ can be chosen, and the algorithm for $S(1)$ is completed since $c_k(1) = 2, \delta k$, and adding any of the remaining messages will violate the constraint. The algorithm then turns to construct $S(2)$ similarly, until all the messages have been chosen, i.e., $E = ?$.

Remark 1: As it can be seen above, the content delivery scheme obtained via Algorithm 1 is not unique. For instance, another feasible content delivery scheme with $s = 2$ for Example 1 is:

$$S(1) = ff1;2g; f3;4g; f1;5g; f2;4g; f4;5gg; \text{ and}$$

$$S(2) = ff1;3g; f2;3g; f3;5g; f1;4g; f2;5gg;$$

Example 2: $N = 4; K = 4; M = 1. t = \frac{MN}{K} = 1$.

We present the content delivery schemes obtained from Algorithm 1 for different values of s :

Case 1: $s = 1$. Algorithm 1 leads to $B = 3$ time slots, each with $\frac{1}{3}n$ channel uses, and

$$S(1) = ff1;2g; f3;4gg;$$

$$S(2) = ff1;3g; f2;4gg;$$

$$S(3) = ff1;4g; f2;3gg;$$

It is noted that the scheme is the same as the one in [16] obtained for $s = 3$ and $t = 1$, in the sense that the same sets of messages are transmitted over the same number of time slots.

Case 2: $s = 2$. Algorithm 1 leads to $B = 2$ time slots, and

$$S(1) = ff1;2g; f3;4gg; ff1;3g; f2;4gg;$$

$$S(2) = ff1;4g; f2;3gg;$$

It is noted that there is no s value to induce a scheme in [16] in this scenario, since s can only take the value of 2 to have $s = 2$, making $t + 1$ not divisible by $t + 1$.

Remark 2: The proposed greedy content delivery scheme can be easily extended by limiting the number of active users as in [16]. Specifically, instead of serving as many users as possible, which is up to K , Algorithm 1 can be applied for a user subset of size $t + 1$ to obtain a content delivery scheme under the constraints on the number of messages to decode. While $\frac{t+1}{t}$ must be an integer to induce a content delivery scheme in [16], Algorithm 1 always provides a delivery scheme for any s value. Therefore, our proposed greedy scheme can be considered as a generalization of the one in [16].

Remark 3: It is noted that Algorithm 1 may lead to unequal number of messages transmitted in different time slots, which can be highly sub-optimal. An intuitive way to enhance the performance is to allocate more channel uses to the time slot with more messages to deliver. In general, once the non-overlapping partition of S , i.e., $\sum_{i=1}^B fS(i)g$, is obtained, we can set the blocklength for the transmission of $S(i)$ proportionally to the number of messages $jS(i)j$. For instance, in the case of $s = 2$ in Example 2, we can allocate $\frac{2n}{3}$ channel uses for $S(1)$ and $\frac{n}{3}$ channel uses for $S(2)$.

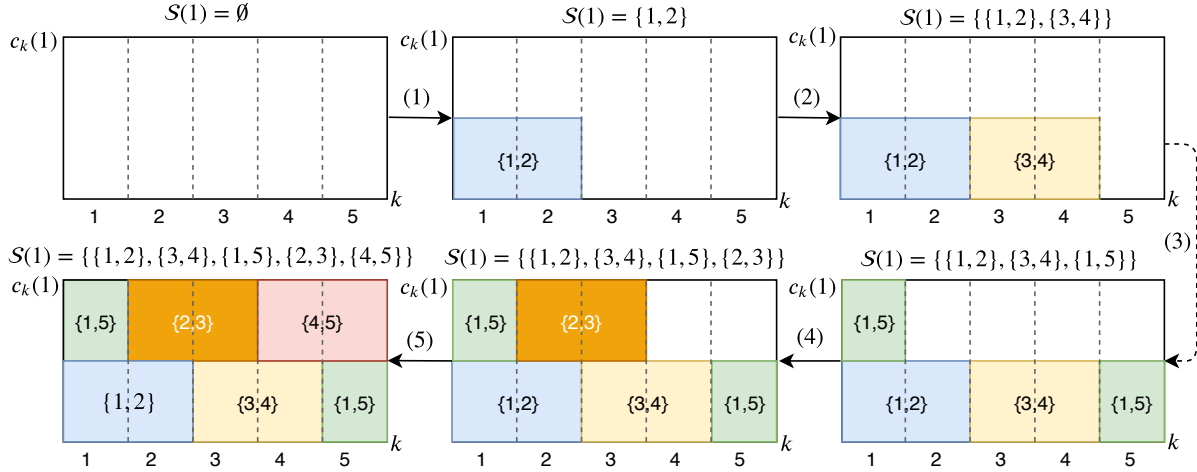


Figure 2: Illustration of the proposed low-complexity greedy scheme for the network with $N = K = 5$, and $M = 1$.

Remark 4: The proposed low-complexity scheme focuses on limiting S , which determines the complexity of multiuser detection at the receivers. On the other hand, the beamformer design resorts to solving the optimization problem in (12), whose computational complexity depends on the content delivery scheme. Given the proposed low-complexity scheme in Algorithm 1, the optimization problem in (12) can be solved as detailed in the Appendix, which is firstly decomposed into at most B_U problems in the form of (32). For each problem, we have at most $K(2^S - 1)$ constraints to characterize the achievable rate region. For the same S value, the optimization problem for the scheme in [16] can be decomposed into $B_I \cdot \frac{K}{t+1} \cdot \frac{(t+1)!}{1!(t+1)!}$ subproblems, each with $(t+1)(2^S - 1)$ constraints for the achievable rate region characterization, where $\frac{t+1}{t+1} = \frac{t+1}{t+1}$. Therefore, each optimization problem in the proposed scheme involves $\frac{K}{t+1}$ times more constraints compared to [16], but may require solving much fewer problems.

Table I compares the proposed scheme in Algorithm 1 and the scheme in [16] in terms of the number of time slots, or equivalently, the number of optimization problems to be solved, denoted by B_U and B_I , respectively, and the ratio of the number of constraints in each of these problems in the former scheme to the latter. We observe that, by simultaneously serving as many users as possible, the optimization problem in each time slot of the proposed low-complexity scheme has comparable number of constraints, but much fewer such optimization problems are needed, as compared to the scheme in [16].

Numerical results for the minimum required power for the proposed greedy transmission scheme, and the comparison with the one proposed in [16] will be presented in Section VII.

V. DEGREES-OF-FREEDOM (DOF) ANALYSIS

In this section we analyze the performance of the scheme proposed in Algorithm 1 in terms of the DoF it achieves in the high SNR regime. To this end, we develop a content delivery scheme which upper bounds the number of time slots B of the scheme presented in Algorithm 1.

For caching factor $t = MK = N$, Algorithm 1 constructs, at each time slot, a set of distinct subsets of users of size $t + 1$, such that no user index appears more than S times in the set. This procedure is repeated until all the $\binom{K}{t+1}$ distinct subsets of users are selected exactly once. While Algorithm 1 aims at minimizing the number of time slots required to deliver all the multicast messages, it is not possible to know in advance how many time slots will be needed. To overcome this uncertainty, we develop a more relaxed scheme, which utilizes a higher number of time slots than the one presented in Algorithm 1.

We assume that the BS is equipped with $N_T = K - t$ antennas, to simultaneously transmit each coded packet of rate $R = \frac{K}{t}$ to the $t + 1$ users that are interested in this message, while zero-forcing it at the remaining $K - t - 1$ users [13]. At each time slot all the K users are targeted, and each user receives no more than S coded packets, and the coded delivery is performed for a total of B time slots. Assuming equal blocklength for different time slots in the high SNR regime, i.e., $n_i = n = B$, $B \geq 2[B]$, since $N_T = K - t$, we can lower bound the per-user DoF as follows [24], [25]:

$$\text{DoF} \geq \frac{K}{SB}; \quad (16)$$

where B is the number of time slots obtained from Algorithm 1.

Next, we present an upper bound on B . We further divide each time slot to S sub-time-slots for the new content delivery approach. At each sub-time-slot, the goal is to create a set of $(t + 1)$ -element subsets of users, such that no user index appears in more than one subset. Due to the symmetry, it is easy to verify that the maximum number of such subsets at each sub-time-slot is $\frac{K}{t+1}$. Therefore, we can generate a set of distinct $(t + 1)$ -element subsets of users, with no user appearing more than S times by repeating this procedure for S sub-time-slots. Accordingly, we can generate a set of $S \cdot \frac{K}{t+1}$ distinct subsets of users, each of size $t + 1$, such that each user index appears no more than S times. This provides us with the set of coded packets for delivery at each time slot. For example, considering Example 1, where $N = 5$; $K =$

Table I: Comparison between Algorithm 1 and [16] in terms of the number of constraints they need to consider for each optimization problem, and the number of time slots, or equivalently, the number of optimization problems to be solved.

$(K; t; s) = (10; 1; 1)$					$(K; t; s) = (10; 2; 1)$					$(K; t; s) = (10; 3; 1)$				
$(;)$	$\frac{K}{t+s}$	B_u	B_l	$\frac{B_u}{B_l}$	$(;)$	$\frac{K}{t+s}$	B_u	B_l	$\frac{B_u}{B_l}$	$(;)$	$\frac{K}{t+s}$	B_u	B_l	$\frac{B_u}{B_l}$
(1;1)	5	9	45	0.2	(1;1)	$\frac{10}{3}$	40	120	0.3333	(1;1)	$\frac{5}{2}$	105	210	0.5
(1;3)	9	9	630	0.0143	(1;4)	2	40	2100	0.0190	(1;5)	$\frac{5}{4}$	105	1575	0.0667
(1;5)	9	9	3150	0.0029	(1;7)	$\frac{10}{7}$	40	2800	0.0143					
(1;7)	9	9	4725	0.0019										
(1;9)	1	9	945	0.0095										

5; $t = 1$, and $s = 2$, we can generate the following sets of subsets of users, each of size 2, in the two sub-time-slots of the first time slot:

$$\begin{aligned} S(1;1) &= f\bar{f}1; 2g; \bar{f}3; 4gg; \text{ and} \\ S(1;2) &= f\bar{f}1; 3g; \bar{f}4; 5gg; \end{aligned} \quad (17)$$

where $S(1) = fS(1;1); S(1;2)g$, and this partitioning is not unique. We exploit the symmetry in the subsets of users of size $t + 1$, to which the coded packets are targeted, to obtain the total number of time slots required for transmission. By creating $s \frac{K}{t+1}$ $(t + 1)$ -element distinct subsets of users at each time slot, we need no more than $\frac{\binom{K}{t+1}}{sB \frac{K}{t+1} C}$ time slots to deliver all the coded messages.

We highlight that, this approach resembles the set generation process in Algorithm 1, but is stricter as it requires each user index to appear no more than once in each sub-time-slot, compared to Line 8 in Algorithm 1. Hence, it results in fewer selected messages at each time slot, and more time slots. Moreover, it is guaranteed in Algorithm 1 that each of the K users decodes at least $s - 1$ messages in each time slot, except for the final time slot. Therefore, the number of time slots is also bounded by $\frac{\binom{K}{t+1}}{s-1} + 1$. Consequently, we have

$$B_u, \min: \frac{82}{6} < \frac{K}{s} \frac{j^{t+1} k}{t+1} \frac{7}{7}; \frac{K-1}{s-1} + 1; \quad (18)$$

and the DoF of the proposed low-complexity scheme satisfies

$$\text{DoF} \frac{K}{sB_u}; \quad (19)$$

The lower bound on the DoF of the proposed greedy scheme is depicted in Fig. 3 for different network parameters, in comparison with the scheme in [16]. As the greedy scheme in Algorithm 1, a low-complexity scheme that achieves the DoF lower bound can always be obtained for any s , while the scheme in [16] may not exist for some values of s , for which the DoF is set to 0 in Fig. 3. Note that the achievable DoF in [16] monotonically increases with s . For any s value, if an integer t is found such that $\frac{t+1}{s} = s$, then we find the maximum possible t such that $t+1$ is divided by $t+1$, yielding the highest DoF of the scheme in [16] as depicted in Fig. 3. It can be seen from Fig. 3 that, the proposed greedy scheme can outperform [16] for certain values of s , especially for small s that is of particular interest in practice. We remark that, the derived lower bound on DoF can be loose due to

the floor and ceiling operations in (18), and the proposed greedy scheme may achieve a DoF strictly higher than the lower bound illustrated here. For instance, in the scenario of $N = K = 9$ and $M = 1$ as considered in Fig. 3(d), the total number of messages for each user to decode is 8. The proposed low-complexity scheme for $s = 8$ delivers all the coded messages in only one time slot, achieving the same DoF of $\frac{9}{8}$ as the scheme in [16]. Hence, the lower bound on DoF at $s = 8$ is very loose as can be seen in Fig. 3(d). We finally remark that, while the DoF lower bound is in general loose, it is tight at certain values of s regardless of the values of $N; K$, and M . Specifically, when $s = 1$, or when s is sufficiently large such that $B_u = 1$, the scheme obtained via the relaxed algorithm, which derives the DoF lower bound in (19), is identical to the one obtained via Algorithm 1 with equal blocklength allocation over all the time slots; and hence, the right-hand-side term in (19) is the exact DoF achieved by the scheme in Algorithm 1 when simply allocating equal blocklength across all the time slots.

VI. JOINT OPTIMIZATION OF BEAMFORMING AND CODED CONTENT DELIVERY

In this section, we formulate a sparsity constrained power minimization problem to jointly optimize the beamformers and the content delivery scheme. The sparsity induced problem directly limits the number of messages to be decoded by each user at any time slot, and the indicator function $v_T(i)$ is identified by setting $v_T(i) = jR^T(i)j_0$; for $8i; T$, where $j j_0$ denotes the ℓ_0 -norm and is equal to the number of non-zero elements of a vector. Therefore, we impose an ℓ_0 -norm constraint on the rates of messages at any time slot i as follows:

$$\times \quad jR^T(i)j_0 \quad s; 8k; i; \quad (20)$$

In this section, we assume equal blocklength allocation over all the B time slots for simplicity. Then, the minimum required power problem with the constraints on the number of messages to be decoded by any user at any time slot can be formulated

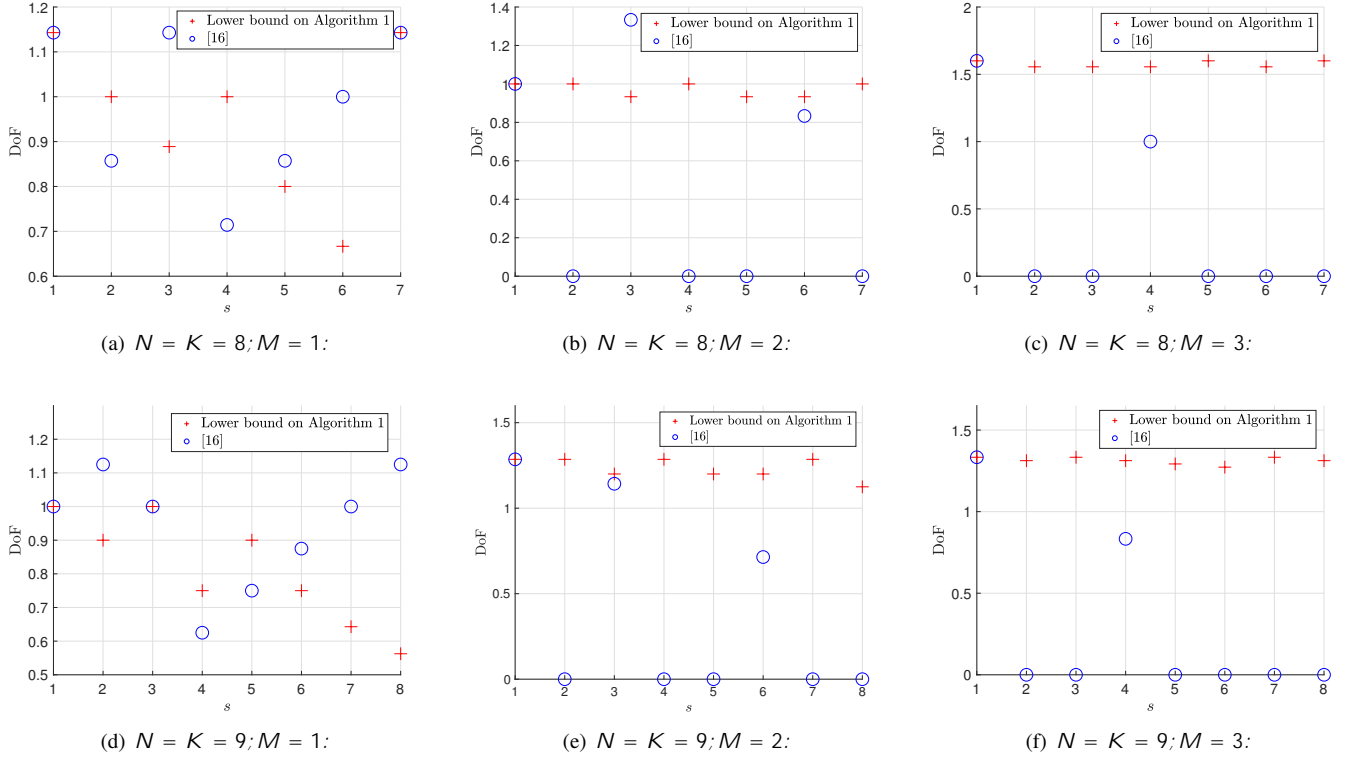


Figure 3: DoF lower bound for the proposed greedy scheme in Algorithm 1 compared to the DoF of the scheme in [16] as a function of s for various network parameters.

as follows:

$$\min_{f_{w_T}(i); f_{R^T}(i); g} \frac{1}{B} \sum_{i=1}^B \sum_{T \in \mathcal{T}_S} k w_T(i) k^2 \quad (21a)$$

$$\text{s.t.} \quad \sum_{T \in \mathcal{T}_{S_k}^j} R^T(i) \leq \frac{1}{B} \log_2 \left(1 + \sum_{T \in \mathcal{T}_{S_k}^j} T_k(i) \right); \quad (21b)$$

$$\sum_{i=1}^B R^T(i) \leq \frac{R}{K}; \quad \forall T; \quad (21c)$$

$$\sum_{T \in \mathcal{T}_{S_k}^j} j R^T(i) \leq s; \quad \forall s; \quad \forall i; \quad (21d)$$

where $T_k(i)$ is defined in (9). In this problem, the objective is to minimize the average transmission power over all the time slots; constraints (21b)-(21c) guarantee successful decoding of all the required messages at each user in each time slot; and constraint (21d) limits the number of messages decoded by each user in any time slot. Since (21d) limits only the number of messages decoded by each user in any time slot, without assuming any specific content delivery scheme, the problem in (21) includes the content delivery schemes in [16] as a special case. This formulation also generalizes the one presented in Section III when the time slots are of equal duration. However, note that the number of time slots B is a free variable for the greedy scheme, while it is assumed to be given in (21).

To deal with the discontinuous ℓ_0 -norm constraint in (21d), we approximate it with a differentiable continuous function [26]

$$f(R^T(i); t) = \frac{2}{\pi} \arctan \frac{R^T(i)}{t}; \quad (22)$$

where $t > 0$ is a prescribed constant that determines the approximation accuracy. The function in (22) is concave w.r.t. $R^T(i)$, therefore the approximate constraint for (21d)

$$\sum_{T \in \mathcal{T}_{S_k}^j} f(R^T(i); t) \leq s; \quad \forall s; \quad \forall i; \quad (23)$$

is concave and can be treated as a difference of convex function. Overall, the original problem in (21) can be approximated by the following problem:

$$\min_{f_{w_T}(i); f_{R^T}(i); g} \frac{1}{B} \sum_{i=1}^B \sum_{T \in \mathcal{T}_S} k w_T(i) k^2 \quad (24a)$$

s.t. (21b); (21c) and (23):

Similarly to (12), the problem in (21) can be solved via the SCA method. Specifically, we introduce $J_{S_k}^j(i)$, then the constraint in (21b) can be rewritten

similarly to (30b) and (30c), given by

$$\times \prod_{T,2}^{j_{S_k}} R^T(i) \frac{1}{B} \log_2(1 + \frac{j_{S_k}(i)}{S_k}; 8k; i; \quad (25)$$

$$\times \prod_{1,2S_k^c} j_{S_k}^H w_l(i) j^2 + \frac{\prod_{T,2}^{j_{S_k}} j_{S_k}^H w_T(i) j^2}{j_{S_k}(i)} + \frac{2}{k} 0; 8k; i; \quad (26)$$

where the constraints in (25) are convex, and the constraints in (26) are in the form of difference of convex functions. In the $(t+1)$ -th iteration of the SCA algorithm, the constraints in (26) can be linearized with the first order Taylor expansion at $f_{w_T}(i)g$ and $f_{j_{S_k}}(i)g$, leading to stricter constraints given by

$$\times \prod_{1,2S_k^c} j_{S_k}^H w_l(i) j^2 + \frac{\prod_{T,2}^{j_{S_k}} j_{S_k}^H w_T(i) j^2}{j_{S_k}(i)} + \frac{2 \prod_{T,2}^{j_{S_k}} w_T^H(i) h_k h_k^H w_T(i)}{j_{S_k}(i)} + \frac{2}{k} 0; 8k; i; \quad (27)$$

where $f_{w_T}(i)g$ and $f_{j_{S_k}}(i)g$ are the solutions to the subproblem in the t -th iteration. The same linearization technique can also be performed for the constraints in (23), yielding stricter constraints given by

$$\times \prod_{T,2S_k} \arctan \frac{R^T(i)}{t} + \frac{t}{t^2 + R^T(i)^2} R^T(i) + \frac{S}{2}; 8k; i; \quad (28)$$

where $f_{R^T}(i)g$ are the solutions to the subproblem in the t -th iteration. Overall, the convex subproblem to be solved in the $(t+1)$ -th iteration is

$$\min_{f_{w_T}(i)g; f_{R^T}(i)g; f_{j_{S_k}}(i)g} \frac{1}{B} \sum_{i=1}^{K} k w_T(i) k^2 \quad (29a)$$

s.t. (21c); (25), (27), and (28):

The initialization of the SCA algorithm for problem (29) for a given value of S requires a content delivery scheme with less or equal complexity, which can be obtained via Algorithm 1 in Section IV. The associated beamforming design can be readily obtained similarly to (31). From problem (21), we can also conclude that the minimum required power of a content delivery scheme is a non-decreasing function of S , as the problem becomes more relaxed as S increases.

VII. SIMULATION RESULTS

We consider a single-cell with radius 500m, and users uniformly randomly distributed in the cell. Channel vectors \mathbf{h}_k are written as $\mathbf{h}_k = (10^{-\text{PL}/10})^{1/2} \mathbf{h}_k$, $8k$, where \mathbf{h}_k denotes an i.i.d. vector accounting for Rayleigh fading of unit power, and the path loss exponent is modeled as $\text{PL} = 148.1 + 37.6 \log_{10}(v_k)$, with v_k denoting the distance between

the BS and the user (in kilometers). The noise variance is set to $\frac{2}{k} = 2 = 134$ dBW for all the users. Throughout this section, we assume that the number of transmit antennas is $N_T = K + t$ as required for achieving a satisfactory DoF performance. All simulation results are averaged over 300 independent trials computed with CVX [27].

The scheme with $B = 1$ time slot will be referred to as the full superposition (FS) scheme. FS has the best performance in terms of transmit power given enough spatial DoF, and serves as a baseline, but it also has the highest complexity. To compare our results with those in [16], same number of coded messages are transmitted to each user in each time slot for both schemes. We note here that with the use of β parameter, the scheme in [16] can be improved by serving disjoint subsets of users simultaneously without increasing the complexity, but the improvement is only applicable when the size of user subset can be partitioned equally and exactly. Therefore, the scheme in [16] cannot handle certain settings such as the case of $s = 2$ in Example 1.

We first present the average transmit power as a function of the target rate R in Fig. 4 for Example 1, assuming that the BS is equipped with $N_T = 6$ antennas. The scheme in [16] that satisfies $s = 2$ is adopted for fair comparison, where $t = 3$ users are served in each time slot. We observe that the proposed greedy scheme provides significant savings in the transmit power compared to [16] at all rates. The power savings increase with rate R as a result of the increased superposition coding gain. Furthermore, the gap between the proposed greedy scheme and FS is quite small, and remains almost constant with rate. At $R = 8$ bps/Hz, the power loss of the scheme in [16] and ours compared to FS are about 8.5 dB and 0.5 dB, respectively. Hence, we can conclude that the proposed greedy scheme provides significant reduction in the computational complexity without sacrificing the performance much.

The average transmit power as a function of file rate R is further investigated for the setting with $N = 6$ files, $K = 6$ users, $M = 1$, and $N_T = 6$ antennas. Similarly to Fig. 4, it is observed in Fig. 5 that our proposed low-complexity greedy scheme substantially outperforms the scheme in [16] with the same value of s in the high SNR regime. For example, the power savings of the greedy delivery scheme compared to [16] are 8dB and 2dB, for $s = 3$ and $s = 4$, respectively, at $R = 10$ bps/Hz. The power gain is again observed to be larger as the rate increases, while in the low SNR/rate regime, all the schemes achieve comparable performance regardless of s . Also, for the proposed greedy scheme, a larger s allows achieving the same rate with lower transmit power in the high SNR regime, at the expense of increased complexity at the receivers. It is noted that the proposed greedy scheme yields the same content delivery scheme in terms of the transmitted coded messages in each time slot as the one in [16] when $s = 1$ and $s = 2$, which correspond to $\beta = 5; \beta = 1$, and $\beta = 5; \beta = 2$ in [16], respectively.

Fig. 6 and Fig. 7 show the average transmit power versus rate R for Example 2, with $N = 4$ files, $K = 4$ users, $M = 1$, and $N_T = 3$ antennas. In Fig. 6, we compare our greedy content delivery scheme with the one obtained

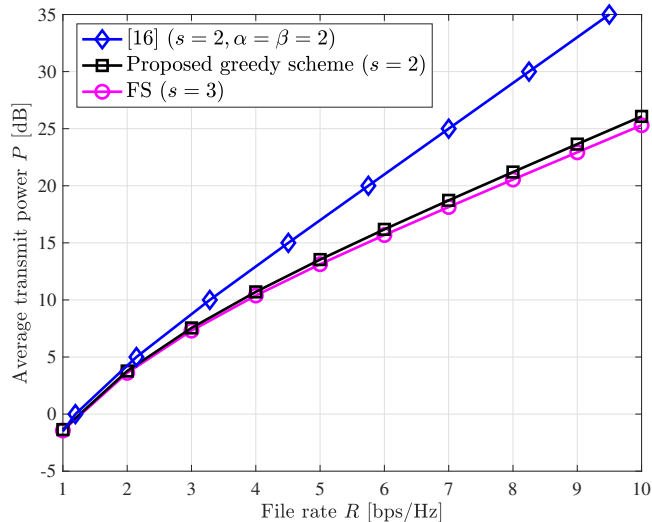


Figure 4: Average transmit power P as a function of rate R for $N = K = 5$, $M = 1$, and $N_T = 6$.

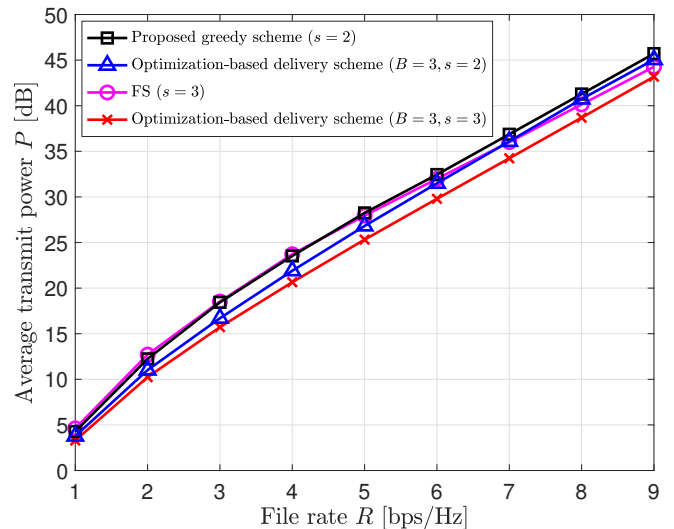


Figure 6: Average transmit power P as a function of rate R for $N = K = 4$, $M = 1$, and $N_T = 3$.

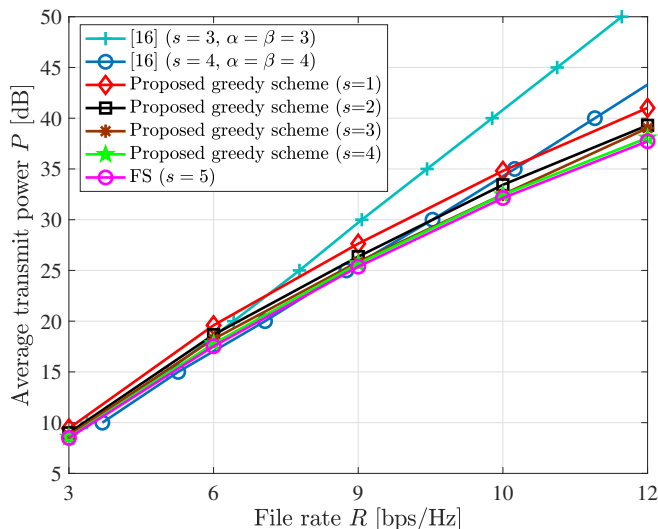


Figure 5: Average transmit power P as a function of rate R for $N = K = 6$, $M = 1$, and $N_T = 6$.

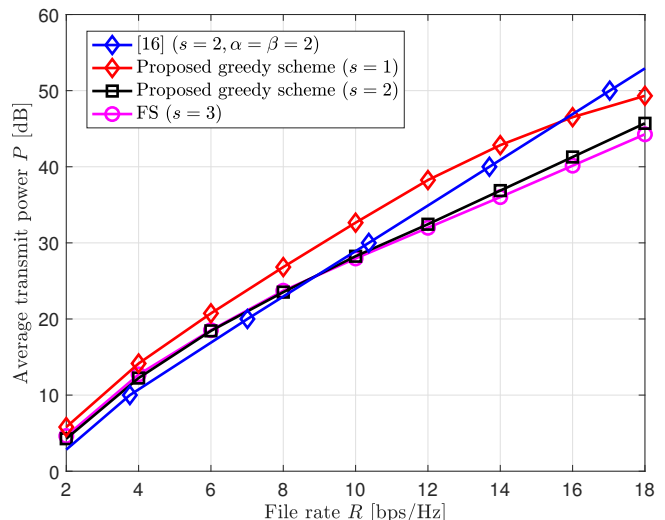


Figure 7: Average transmit power P as a function of rate R for $N = K = 4$, $M = 1$, and $N_T = 3$.

by solving the problem in (21) for $s = 2$ and $s = 3$, by setting $B = 3$. It is seen that the greedy scheme can achieve comparable performance, and the performance gap is small especially for high rates. The optimization-based content delivery scheme with $s = 3$ is found to outperform the one with $s = 2$ as expected, and the improvement is larger as the rate increases. In Fig. 7, it is interesting to see that when the rate is low, the scheme in [16] slightly outperforms both the FS and the proposed schemes. A similar observation has been made in [16] showing that a higher rate can be achieved when transmitting a smaller number of coded messages at low SNR, reducing the interference from coded packets that are not decoded at each user. Due to insufficient spatial degrees of freedom, both the FS and the proposed schemes fail to manage

the interference between data streams. We conclude that this effect occurs only for low rates, as the benefit of superposition coding becomes more dominant at higher rates. We note that the greedy scheme coincides with the one in [16] for $s = 1$ in terms of the transmitted coded messages in each time slot, but this does not always happen. For instance, when $s = 2$, the only option in [16] to keep the same level of complexity is to serve 3 users in each time slot.

We plot in Fig. 8 the power loss of the proposed scheme in Algorithm 1 compared to FS as a function of s , that is, how much more the required transmit power is compared to the transmit power required by FS. This results in a more clear plot. Assuming $N = 6$, $K = 6$, $M = 1$, we let s take values from $\{1; 2; 3; 4; 5\}g$, where $s = 5$ corresponds to the

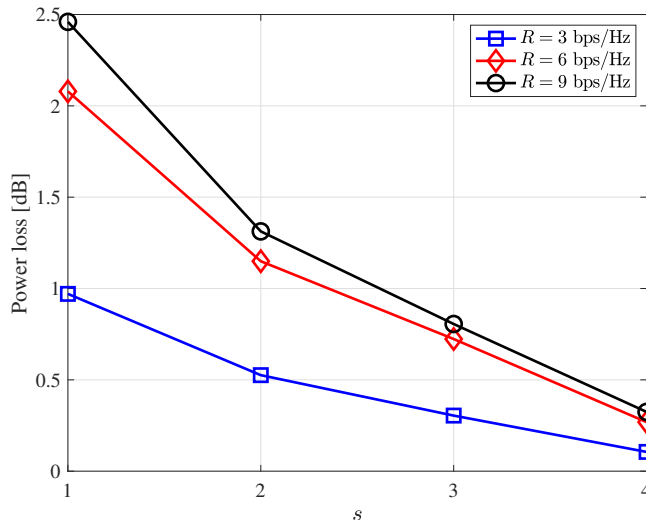


Figure 8: Power loss w.r.t. FS as a function of s for $N = K = 6$, $M = 1$, and $N_T = 6$.

FS scheme in which all the $\binom{K}{t+1} = 15$ coded messages are transmitted simultaneously. FS requires less transmit power at the expense of a high decoding complexity. When $s = 1$, the model boils down to the single-cell multigroup multicasting problem, which has the lowest computation and implementation complexity. In general, Fig. 8 can be considered as the trade-off curve between the performance and complexity for each rate value, both of them increasing with s .

VIII. CONCLUSIONS

In this paper, we have studied cache-aided content delivery from a multi-antenna BS in the finite SNR regime. We have formulated a general beamforming scheme that multicasts coded files over multiple orthogonal time slots. We have then specialized this general formulation to a low-complexity greedy scheme by limiting the number of coded messages targeted at each user at each time slot. This scheme provides the flexibility to adjust the computational complexity of the optimization problem and the receiver complexity. We have then formulated the constraint on the number of coded messages targeted at each user at each time slot as a sparsity constraint, and solved the resulting mixed-integer non-convex optimization problem using the SCA method. Compared with FS, where all the coded messages are transmitted simultaneously, and the scheme obtained via the sparsity-constrained optimization framework, the greedy scheme achieves comparable performance, and outperforms the one proposed in [16] for all values of SNR and rate with sufficient spatial degrees of freedom, while the improvement is limited to high data rate values when the BS does not have sufficiently many transmit antennas. Furthermore, the gap between the greedy delivery scheme and the optimization-based delivery scheme decreases as the SNR/power increases. When considering practical implementations, one must choose a suitable value of s that yields an acceptable performance while keeping the complexity feasible. The satisfactory DoF performance of the

proposed low-complexity scheme is guaranteed with at least $K - t$ antennas, while the analysis of overloaded systems with K users served by less than $K - t$ antennas is left as future work.

APPENDIX A

AN UPPER BOUND ON THE SOLUTION OF PROBLEM (12)

It is noted that the constraints are in the form of difference of convex functions, which can be approximated by linearizing the concave functions, resulting in a convex problem that can be solved via SCA techniques. To see this, we first rewrite the problem in (12) as

$$\min_{\mathbf{w}_T(i); \mathbf{R}^T(i); \mathbf{f}_{S_k}^{(i)}; g_{i=1}^{T, S}} \times \prod_{i=1}^{T, S} \frac{n_i}{n} k \mathbf{w}_T(i) k^2 \quad (30a)$$

$$\text{s.t.} \quad \times \prod_{T, S} \mathbf{R}^T(i) \frac{n_i}{n} \log_2(1 + \frac{j_{S_k}(i)}{j_{S_k}});$$

$$\times \prod_{T, S} \frac{8 \sum_{S_k}^2 \mathbf{h}_k^H \mathbf{w}_T(i) \mathbf{h}_k}{j_{S_k}^2} + \frac{2}{k} 0; 8 \sum_{S_k}^2 \mathbf{h}_k^H \mathbf{w}_T(i) \mathbf{h}_k; i; \quad (30b)$$

$$\times \prod_{T, S} \frac{j_{S_k}^H \mathbf{w}_T(i) \mathbf{h}_k^H \mathbf{w}_T(i) \mathbf{h}_k}{j_{S_k}^2} + \frac{2}{k} 0; 8 \sum_{S_k}^2 \mathbf{h}_k^H \mathbf{w}_T(i) \mathbf{h}_k; i; \quad (30c)$$

$$(12c), (12d) \text{ and } (12e);$$

where $\frac{j_{S_k}(i)}{j_{S_k}}$, $\mathbf{R}^T(i)$. The constraint in (30c) is the difference of convex function, since $\prod_{T, S} \frac{j_{S_k}^H \mathbf{w}_T(i) \mathbf{h}_k^H \mathbf{w}_T(i) \mathbf{h}_k}{j_{S_k}^2}$ and $\frac{j_{S_k}(i)}{j_{S_k}}$ is the sum of quadratic-over-linear functions of $\mathbf{w}_T(i)$ and $\frac{j_{S_k}(i)}{j_{S_k}}$. Therefore, a sequence of convex subproblems can be solved iteratively to approximately tackle this convex-concave problem [28], with the subproblem in the $(+1)$ -th iteration given by

$$\min_{\mathbf{w}_T(i); \mathbf{R}^T(i); \mathbf{f}_{S_k}^{(i)}; g_{i=1}^{T, S}} \times \prod_{i=1}^{T, S} \frac{n_i}{n} k \mathbf{w}_T(i) k^2 \quad (31a)$$

$$\text{s.t.} \quad \times \prod_{T, S} \mathbf{R}^T(i) \frac{n_i}{n} \log_2(1 + \frac{j_{S_k}(i)}{j_{S_k}}); 8 \sum_{S_k}^2 \mathbf{h}_k^H \mathbf{w}_T(i) \mathbf{h}_k; i; \quad (31b)$$

$$\times \prod_{T, S} \frac{j_{S_k}^H \mathbf{w}_T(i) \mathbf{h}_k^H \mathbf{w}_T(i) \mathbf{h}_k}{j_{S_k}^2} + \frac{2}{k} 0; 8 \sum_{S_k}^2 \mathbf{h}_k^H \mathbf{w}_T(i) \mathbf{h}_k; i; \quad (31c)$$

(12c), (12d) and (12e);

given the solution of $\mathbf{w}_T(i)$, $\mathbf{R}^T(i)$, and $\frac{j_{S_k}(i)}{j_{S_k}}$ obtained in the $(-)$ -th SCA iteration. Each of the convex subproblems can be efficiently solved with standard interior-point algorithms or off-the-shelf solvers, and the SCA approach is guaranteed to converge to a stationary solution of the original problem in (12) [29]. Details of the SCA algorithm are outlined in Table. II.

An initial point in the feasible set of problem (12) is required to initialize the SCA algorithm. We first observe that for any feasible target rates $fR^T(i)j\delta T \geq Sg_{i=1}^B$ that satisfy the constraints in (12c) and (12e), the problem in (12) can be decoupled and decomposed into B parallel subproblems, each for a distinct time slot $i \in [B]$, given by

$$fW_T(i)g_{T \geq S(i)} = \arg \min_{fW_T(i)g_{T \geq S(i)}} \sum_{T \geq S(i)} kW_T(i)k^2 \quad (32a)$$

$$\text{s.t.} \quad \sum_{T \geq S_k} R^T(i) \leq \frac{n_i}{n} \log_2 \left(1 + \sum_{T \geq S_k} T_k(i) \right); \quad \sum_{T \geq S_k} T_k(i) \leq S_k; \quad \delta k; i; \quad (32b)$$

$$T_k(i) = \mathbb{P} \frac{j\mathbf{h}_k^H W_T(i)j^2}{\sum_{l \in S_k} j\mathbf{h}_k^H W_l(i)j^2 + \frac{1}{k}}; \quad \delta k; i; \quad (32c)$$

$$kW_T(i)k^2 = 0 \text{ for } \delta V_T(i) \neq 0; \quad (32d)$$

$$kW_T(i)k^2 = 0 \text{ for } \delta V_T(i) = 0; \quad (32e)$$

which is nonconvex. Nevertheless, it can be transformed into a semidefinite programming problem by introducing $W_T(i)$, $W_T(i)W_T^H(i)$ and dropping the rank-1 constraints on $W_T(i)$, which is given by

$$fW_T(i)g_{T \geq S(i)} = \arg \min_{fW_T(i)g_{T \geq S(i)}} \text{Tr}fW_T(i)g \quad (33a)$$

$$\text{s.t.} \quad \sum_{T \geq S_k} R^T(i) \leq \frac{n_i}{n} \log_2 \left(1 + \sum_{l \in S_k} \text{Tr}f\mathbf{H}_k W_l(i)g + \frac{1}{k} \right); \quad \sum_{T \geq S_k} \text{Tr}f\mathbf{H}_k W_T(i)g \leq S_k; \quad \delta k; i; \quad (33b)$$

$$W_T(i) = 0; \delta V_T(i) \neq 0; \quad (33c)$$

$$W_T(i) = 0; \delta V_T(i) = 0; \quad (33d)$$

and can be efficiently solved with standard interior-point algorithms. However, the solution obtained with semidefinite relaxation is not necessarily rank-1. If the obtained $W_T(i)$'s are all rank-1, then the optimal solution of (32) can be readily recovered from $W_T(i)$. Otherwise, Gaussian randomization can be adopted to obtain a feasible approximation to the optimal solution of (32). Note that the solution given by (32) is an upper bound on the minimum required power in (12) as the rates $fR^T(i)j\delta T \geq Sg_{i=1}^B$ are not optimized, which hence can serve as an initial point in the successive convex approximation algorithm to obtain a tighter upper bound on the problem in (12).

REFERENCES

- [1] J. Zhao, M. M. Amiri, and D. Gündüz, "A low-complexity cache-aided multi-antenna content delivery scheme," in *2019 IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 2019, pp. 1–5.
- [2] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. Inf. Theory*, vol. 60, no. 5, pp. 2856–2867, May. 2014.
- [3] S. Saeedi Bidokhti, M. Wigger, and R. Timo, "Noisy broadcast networks with receiver caching," *IEEE Trans. Inf. Theory*, vol. 64, no. 11, pp. 6996–7016, Nov. 2018.

Table II: SCA Algorithm for the Multicast Beamforming Problem with a Given Coded Delivery Scheme

STEP 0: Set $\epsilon = 1$. Set a step size ϵ . Initialize $W_T(i)$, $R_T(i)$, and $\sum_{S_k} j(i)$ with feasible values
STEP 1: If a stopping criterion is satisfied, then STOP
STEP 2: Solve the optimization problem in (31)
STEP 3: Update $W_T^{+1}(i) = W_T(i) + \epsilon W_T(i)$, $W_T(i)$, $R_T^{+1}(i) = R_T(i) + \epsilon R_T(i)$, $R_T(i)$ \downarrow , $\sum_{S_k} j^{+1}(i) = \sum_{S_k} j(i) + \epsilon \sum_{S_k} j(i)$, $\sum_{S_k} j(i)$, ϵ
STEP 4: Set $\epsilon = \epsilon + 1$, and go to STEP 1

- [4] M. Mohammadi Amiri and D. Gündüz, "Cache-aided content delivery over erasure broadcast channels," *IEEE Trans. Commun.*, vol. 66, no. 1, pp. 370–381, Jan. 2018.
- [5] J. Zhang and P. Elia, "Wireless coded caching: A topological perspective," *CoRR*, vol. abs/1606.08253, 2016. [Online]. Available: <http://arxiv.org/abs/1606.08253>
- [6] M. Mohammadi Amiri and D. Gündüz, "Caching and coded delivery over Gaussian broadcast channels for energy efficiency," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 8, pp. 1706–1720, Aug. 2018.
- [7] A. Ghorbel, K. Ngo, R. Combes, M. Kobayashi, and S. Yang, "Opportunistic content delivery in fading broadcast channels," in *GLOBECOM 2017 - 2017 IEEE Global Communications Conference*, Dec. 2017, pp. 1–6.
- [8] N. D. Sidiropoulos, T. N. Davidson, and Z.-Q. Luo, "Transmit beamforming for physical-layer multicasting," *IEEE Trans. Signal Process.*, vol. 54, no. 6, pp. 2239–2251, Jun. 2006.
- [9] E. Karipidis, N. D. Sidiropoulos, and Z. Luo, "Quality of service and max-min fair transmit beamforming to multiple cochannel multicast groups," *IEEE Trans. Signal Process.*, vol. 56, no. 3, pp. 1268–1279, Mar. 2008.
- [10] Z. Xiang, M. Tao, and X. Wang, "Coordinated multicast beamforming in multicell networks," *IEEE Trans. Wireless Commun.*, vol. 12, no. 1, pp. 12–21, Jan. 2013.
- [11] K. Ngo, S. Yang, and M. Kobayashi, "Scalable content delivery with coded caching in multi-antenna fading channels," *IEEE Trans. Wireless Commun.*, vol. 17, no. 1, pp. 548–562, Jan. 2018.
- [12] S. P. Shariatpanahi, G. Caire, and B. Hossein Khalaj, "Physical-layer schemes for wireless coded caching," *IEEE Trans. Inf. Theory*, vol. 65, no. 5, pp. 2792–2807, May. 2019.
- [13] S. P. Shariatpanahi, S. A. Motahari, and B. H. Khalaj, "Multi-server coded caching," *IEEE Trans. Inf. Theory*, vol. 62, no. 12, pp. 7253–7271, Dec. 2016.
- [14] E. Lampiris and P. Elia, "Adding transmitters dramatically boosts coded-caching gains for finite file sizes," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 6, pp. 1176–1188, Jun. 2018.
- [15] M. Salehi, A. Tölli, and S. P. Shariatpanahi, "A multi-antenna coded caching scheme with linear subpacketization," *CoRR*, vol. abs/1910.10384, 2019. [Online]. Available: <http://arxiv.org/abs/1910.10384>
- [16] A. Tölli, S. P. Shariatpanahi, J. Kaleva, and B. H. Khalaj, "Multi-antenna interference management for coded caching," *IEEE Trans. Wireless Commun.*, vol. 19, no. 3, pp. 2091–2106, Mar. 2020.
- [17] E. Piovano, H. Joudeh, and B. Clerckx, "On coded caching in the overloaded MISO broadcast channel," in *2017 IEEE International Symposium on Information Theory (ISIT)*, Jun. 2017, pp. 2795–2799.
- [18] J. Zhang and P. Elia, "Fundamental limits of cache-aided wireless BC: Interplay of coded-caching and CSIT feedback," *IEEE Trans. Inf. Theory*, vol. 63, no. 5, pp. 3142–3160, May. 2017.
- [19] S. Zhong and X. Wang, "Joint multicast and unicast beamforming for coded caching," *IEEE Trans. Commun.*, vol. 66, no. 8, pp. 3354–3367, Aug. 2018.
- [20] T. X. Vu, S. Chatzinotas, and B. Ottersten, "Edge-caching wireless networks: Performance analysis and optimization," *IEEE Trans. Wireless Commun.*, vol. 17, no. 4, pp. 2827–2839, Apr. 2018.
- [21] T. X. Vu, L. Lei, S. Chatzinotas, B. Ottersten, and T. A. Vu, "Energy efficient design for coded caching delivery phase," in *2019 3rd International Conference on Recent Advances in Signal Processing, Telecommunications Computing (SigTelCom)*, Mar. 2019, pp. 165–169.
- [22] B. R. Marks and G. P. Wright, "A general inner approximation algorithm

