

# Eigen-Based Transceivers for the MIMO Broadcast Channel with Semi-Orthogonal User Selection

Liang Sun, *Student Member, IEEE* and Matthew R. McKay, *Member, IEEE*

## Abstract

This paper studies the sum rate performance of two low complexity eigenmode-based transmission techniques for the MIMO broadcast channel, employing greedy semi-orthogonal user selection (SUS). The first approach, termed ZFDPC-SUS, is based on zero-forcing dirty paper coding; the second approach, termed ZFBF-SUS, is based on zero-forcing beamforming. We first employ new analytical methods to prove that as the number of users  $K$  grows large, the ZFDPC-SUS approach can achieve the optimal sum rate scaling of the MIMO broadcast channel. We also prove that the average sum rates of both techniques converge to the average sum capacity of the MIMO broadcast channel for large  $K$ . In addition to the asymptotic analysis, we investigate the sum rates achieved by ZFDPC-SUS and ZFBF-SUS for finite  $K$ , and show that ZFDPC-SUS has significant performance advantages. Our results also provide key insights into the benefit of multiple receive antennas, and the effect of the SUS algorithm. In particular, we show that whilst multiple receive antennas only improves the asymptotic sum rate scaling via the second-order behavior of the multi-user diversity gain; for finite  $K$ , the benefit can be very significant. We also show the interesting result that the semi-orthogonality constraint imposed by SUS, whilst facilitating a very low complexity user selection procedure, asymptotically does not reduce the multi-user diversity gain in either first ( $\log K$ ) or second-order ( $\log \log K$ ) terms.

Copyright (c) 2010 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to [pubs-permissions@ieee.org](mailto:pubs-permissions@ieee.org).

Manuscript received Oct. 27, 2009; revised Feb. 10, 2010. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Ali Ghrayeb. L. Sun and M. R. McKay are with the ECE Department, Hong Kong University of Science and Technology, Hong Kong. (Email: [sunliang@ust.hk](mailto:sunliang@ust.hk); [eemckay@ust.hk](mailto:eemckay@ust.hk)). The work of L. Sun and M. R. McKay was supported by the Hong Kong Research Grants Council (RGC) under grant no. 617108. This work was presented in part at the IEEE Global Communications Conference (Globecom), Honolulu, USA, December, 2009.

## I. INTRODUCTION

In the multiple-input multiple-output (MIMO) broadcast channel, the spatial multiplexing capability of multiple transmit antennas can be exploited to efficiently serve multiple users simultaneously, rather than trying to maximize the capacity of a single-user link. The capacity region of the MIMO broadcast channel has now been well-studied [1–5], and has been shown to be achieved through the use of multiple antenna dirty paper coding (DPC) [3]. Unfortunately, optimal DPC is a highly non-linear technique involving joint optimization over a set of power-constrained covariance matrices, and is therefore too complex for practical implementation [4]. A reduced complexity sub-optimal DPC scheme, known as zero-forcing dirty paper coding (ZFDPC), was proposed for single-antenna users in [5], and generalized to multiple-antenna users in [6], which is based on a QR decomposition of the channel matrix.

To further reduce complexity, linear processing schemes such as beamforming (BF) have also attracted a lot of attention. The zero-forcing beamforming (ZFBF) scheme was first introduced for single-antenna users in [5], and further modified in [7] and [8]. In [9], the concept of block-diagonalization was proposed for multiple-antenna users, which completely cancels the inter-user interference by employing a set of precoding matrices. One key limitation of these techniques is that, for ZFDPC and ZFBF, the maximum number of users that can be supported must be no more than the number of transmit antennas, whereas for block-diagonalization, the number of the transmit antennas must be larger than the aggregate number of receive antennas across all users. This is significant, since the number of users in practice can be large.

When the number of users  $K$  is larger than the number of transmit antennas  $M$ , one must select a subset of users in the system. A common approach is to seek the subset of users which yields the maximum sum rate. The complexity of finding the optimal subset, however, can be prohibitively large, and to reduce complexity greedy algorithms are commonly employed (see e.g., [10–12]). A promising way to further reduce the complexity of user selection is to restrict the searching space of users by imposing some constraint on the channels of the selected users. Following this method, [13] proposed a semi-orthogonal user selection (SUS) algorithm which iteratively searches for users with nearly orthogonal channel directions<sup>1</sup>.

In this paper, we consider low complexity transmission and user selection techniques for the MIMO broadcast channel with multiple-antenna users. It is still not clear how much advantage can be gained by employing multiple-antennas at the user terminals. Some recent exceptions which deal with the multiple-antenna user scenario are presented in [14] and [15]. Particularly, [14] proposed a generalized G-ZFDPC approach, based on the idea of eigenmode transmission (eigen-beamforming). A limitation of

<sup>1</sup>More specifically, two complex vectors  $\mathbf{u}$  and  $\mathbf{v}$ , with unit norm, are said to be semi-orthogonal if  $|\mathbf{u}^H \mathbf{v}|^2 < \delta$ , where  $\delta$  is referred to as the *semi-orthogonality parameter*.

that approach is the relatively high complexity, since it requires numerical optimization of certain system parameters. In [15], a thresholding technique based on the channel singular values was proposed, and necessary and sufficient conditions were given to achieve the optimum sum capacity of DPC as  $K \rightarrow \infty$ . However, for that scheme, the optimal threshold must be computed by exhaustive search, and is once again quite complicated when the number of users is not small.

In this paper, we investigate two low complexity eigen-beamforming-based transceiver structures for the MIMO broadcast channel with multiple-antenna users, combined with a greedy SUS algorithm. The first technique is a generalization the G-ZFDPC approach in [10] to account for multiple-antenna users and combine it with SUS. We refer to this technique as ZFDPC-SUS. The second technique is a generalization of the algorithm proposed in [13], which we refer to as ZFBF-SUS. For both techniques, we present an asymptotic performance analysis of the sum rate (as in [6, 13–17]) as the number of users grows large. In particular, by employing novel analytical techniques, we demonstrate that ZFDPC-SUS achieves the optimal sum capacity scaling of the MIMO broadcast channel as the number of users grows large. In addition, we prove the more powerful result that the difference between the sum rate of ZFDPC-SUS and the sum capacity of the MIMO broadcast channel converges to zero. We also establish a similar result for ZFBF-SUS. In addition to the asymptotic analysis, we also investigate the sum rates achieved by ZFDPC-SUS and ZFBF-SUS for finite  $K$ , for high and low signal-to-noise ratios (SNR). Based on our analytical results, we establish a number of important insights. For example, we demonstrate that by employing multiple-antennas at the user terminals only affects the asymptotic sum rate scaling via the second-order behavior of the multi-user diversity gain. Thus, the improvement due to having multiple receive antennas at the terminals is much less than that of having multiple transmit antennas, which provides linear capacity growth through spatial multiplexing gain. However, for finite  $K$ , we show that the performance improvement due to multiple receive antennas can still be very significant. We also establish key insights into the design of the semi-orthogonality parameter used in the SUS algorithm. In particular, it has been claimed previously that the semi-orthogonality constraint will cause multi-user diversity gain reduction [13]. However, through our asymptotic analysis, we show that if some very mild conditions on the semi-orthogonality constraint are met, then the semi-orthogonality parameter *does not* reduce the multi-user diversity gain in either first or second order, for both ZFDPC-SUS and ZFBF-SUS. It seems that this conclusion cannot be established by using previous analytical methods for SUS [13]. Our analysis also leads to practical design guidelines for selecting the semi-orthogonality parameter for finite numbers of users, in order to intelligently trade off complexity and performance. Our analysis also demonstrates that for finite values of  $K$ , ZFDPC-SUS can significantly outperform ZFBF-SUS.

## II. CHANNEL AND SYSTEM MODEL

We consider a MIMO broadcast channel with  $M$  transmit antennas and  $K$  users, with  $K \geq M$ . User  $k$  is equipped with  $N_k$  antennas. In a flat-fading environment, the baseband model of this system is

$$\mathbf{y}_k = \mathbf{H}_k \mathbf{s} + \mathbf{n}_k, \quad 1 \leq k \leq K, \quad (1)$$

where  $\mathbf{y}_k \in \mathcal{C}^{N_k \times 1}$  is the received signal vector of user  $k$ ,  $\mathbf{H}_k \in \mathcal{C}^{N_k \times M}$  denotes the channel matrix from the transmitter to user  $k$ ,  $\mathbf{s} \in \mathcal{C}^{M \times 1}$  represents the transmit signal vector, designed to meet the total power constraint  $\text{Tr}(\mathcal{E}\{\mathbf{s}\mathbf{s}^H\}) \leq P$ , and  $\mathbf{n}_k \in \mathcal{C}^{N_k \times 1}$  is white Gaussian noise with zero mean and covariance matrix  $\mathbf{I}_{N_k}$ . Throughout the paper, we assume (as in [5, 13, 14, 18]) that (i) the channels of all users are subject to uncorrelated Rayleigh fading and, for simplicity, all users are homogeneous and experience statistically independent fading, (ii) the transmitter has perfect CSI of all downlink channels<sup>2</sup>, and (iii) each user only has access to their own CSI, but not the CSI of the downlink channels of the other users.

The transmitter supports  $L \leq M$  simultaneous data streams, shared by at most  $L$  selected users (active users), which are indexed by  $\pi(i)$ ,  $i = 1, 2, \dots, L$ . (Note that the specific user selection algorithm will be discussed in Section III.) The transmitted signal vector is represented as

$$\mathbf{s} = \mathbf{W}\mathbf{P}^{\frac{1}{2}}\mathbf{x}, \quad (2)$$

where  $\mathbf{x} = [x_1, x_2, \dots, x_L]^T$  collects the zero-mean circularly symmetric complex Gaussian information signals for each of the  $L$  data streams, satisfying  $\mathcal{E}\{\mathbf{x}\mathbf{x}^H\} = \mathbf{I}_L$ ,  $\mathbf{P} = \text{diag}\{p_1, p_2, \dots, p_L\}$  accounts for the power loading across the multiple streams, chosen to satisfy  $\sum_{i=1}^L p_i \leq P$ , and  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_L] \in \mathcal{C}^{M \times L}$  represents the precoder matrix, with  $\mathbf{w}_i$  denoting the beamforming vector for the  $i$ -th stream (i.e. for user  $\pi(i)$ ), normalized to satisfy  $\|\mathbf{w}_i\|^2 = 1$ . Note that with this formulation, a given user may be assigned multiple data streams.

From (2), the received signal vector for user  $k$  can be rewritten as

$$\mathbf{y}_k = \mathbf{H}_k \mathbf{W}\mathbf{P}^{\frac{1}{2}}\mathbf{x} + \mathbf{n}_k. \quad (3)$$

It is convenient to represent  $\mathbf{H}_k$  via its singular value decomposition (SVD)  $\mathbf{H}_k = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^H$ , where  $\mathbf{\Sigma}_k$  is a  $N_k \times M$  diagonal matrix containing the singular values of  $\mathbf{H}_k$  in decreasing order along its main diagonal, and  $\mathbf{U}_k = [\mathbf{u}_{k,1}, \mathbf{u}_{k,2}, \dots, \mathbf{u}_{k,N_k}] \in \mathcal{C}^{N_k \times N_k}$  and  $\mathbf{V}_k = [\mathbf{v}_{k,1}, \mathbf{v}_{k,1}, \dots, \mathbf{v}_{k,M}] \in \mathcal{C}^{M \times M}$  are

<sup>2</sup>This assumption is reasonable in time division duplex (TDD) systems, which allows the transmitter to employ reciprocity to estimate the downlink channels.

unitary matrices with  $\mathbf{u}_{k,j}$  and  $\mathbf{v}_{k,j}$  representing the left and right singular vectors corresponding to the  $j$ -th largest singular value  $\sqrt{\lambda_{k,j}}$ .

To detect the data stream  $i$ , user  $\pi(i)$  left multiplies the received vector by  $\mathbf{u}_{\pi(i),d_i}$  as follows

$$\begin{aligned} r_{\pi(i),d_i} &= \mathbf{u}_{\pi(i),d_i}^H \mathbf{y}_{\pi(i)} \\ &= \sqrt{\lambda_{\pi(i),d_i}} \mathbf{v}_{\pi(i),d_i}^H \mathbf{W} \mathbf{P}^{\frac{1}{2}} \mathbf{x} + \tilde{n}_{\pi(i),d_i}, \end{aligned} \quad (4)$$

where  $\tilde{n}_{\pi(i),d_i} = \mathbf{u}_{\pi(i),d_i}^H \mathbf{n}_{\pi(i)} \sim \mathcal{CN}(0, 1)$  is the effective additive white Gaussian noise after processing, and  $d_i$  denotes the *eigen-mode index* for stream  $i$ , chosen according to the selection procedure outlined in Section III. Collecting the processed signals (4) for each of the  $L$  data streams, we may write

$$\mathbf{r} = \mathbf{C}_{\pi,d} \mathbf{W} \mathbf{P}^{\frac{1}{2}} \mathbf{x} + \tilde{\mathbf{n}} = \mathbf{\Lambda}_{\pi,d}^{\frac{1}{2}} \mathbf{\Xi}_{\pi,d} \mathbf{W} \mathbf{P}^{\frac{1}{2}} \mathbf{x} + \tilde{\mathbf{n}}, \quad (5)$$

where  $\mathbf{C}_{\pi,d} = [\mathbf{c}_{\pi(1),d_1}^T, \mathbf{c}_{\pi(2),d_2}^T, \dots, \mathbf{c}_{\pi(L),d_L}^T]^T$  is the composite channel matrix for the selected users and eigen-channel set with  $i$ -th row vector  $\mathbf{c}_{\pi(i),d_i} = \sqrt{\lambda_{\pi(i),d_i}} \mathbf{v}_{\pi(i),d_i}^H$ ,  $\tilde{\mathbf{n}} = [\tilde{n}_{\pi(1),d_1}, \tilde{n}_{\pi(2),d_2}, \dots, \tilde{n}_{\pi(L),d_L}]^T$ ,  $\mathbf{\Lambda}_{\pi,d} = \text{diag}\{\lambda_{\pi(1),d_1}, \dots, \lambda_{\pi(L),d_L}\}$ , and  $\mathbf{\Xi}_{\pi,d} = [\mathbf{v}_{\pi(1),d_1}, \dots, \mathbf{v}_{\pi(L),d_L}]^H$ .

In the next section, we will describe several transceiver structures, as well as a greedy method for selecting the set of active users  $\pi = \{\pi(1), \dots, \pi(L)\}$  and the corresponding eigen-channels (active eigen-channels)  $d = \{d_1, \dots, d_L\}$ .

### III. TRANSCEIVER STRUCTURES AND USER SELECTION ALGORITHM

#### A. Greedy Zero-Forcing Dirty Paper Coding Algorithm

In this subsection, we present a transmission strategy which jointly combines ZF, DPC, and eigen-beamforming, along with a greedy low complexity SUS scheduling algorithm. Henceforth, this strategy will be termed ZFDPC-SUS. To the best of our knowledge this scheme has not been considered before. We note, however, that it is an extension of the ZFDPC strategy considered in [5, 10, 18] to account for multiple receive antennas, and also a variation of the algorithm discussed briefly in [13, Sect. VIII].

Let  $\mathbf{\Xi}_{\pi,d} = \mathbf{L}_{\pi,d} \mathbf{Q}_{\pi,d}$  denote the QR decomposition of  $\mathbf{\Xi}_{\pi,d}$ , where  $\mathbf{L}_{\pi,d}$  is a  $L \times L$  lower triangular matrix with  $(i, j)$ -th entry  $l_{i,j}$ , and  $\mathbf{Q}_{\pi,d} = [\mathbf{q}_1^T, \dots, \mathbf{q}_L^T]^T$  is a  $L \times M$  matrix with orthonormal rows ( $\mathbf{q}_i$  denotes the  $i$ -th row vector). The transmit precoder matrix is chosen as

$$\mathbf{W} = \mathbf{Q}_{\pi,d}^H. \quad (6)$$

Then, (5) yields a set of interference channels

$$r_{\pi(i),d_i} = \sqrt{\lambda_{\pi(i),d_i}} (\sqrt{p_i} l_{i,i} x_i + \sum_{j < i} \sqrt{p_j} l_{i,j} x_j) + \tilde{n}_{\pi(i),d_i}. \quad (7)$$

From (7), if  $i < j$ , there is no interference at receiver  $\pi(i)$  from data stream  $j$ . For  $i > j$ , the interference term  $\sum_{j < i} \sqrt{p_j} l_{i,j} x_j$  is precanceled at the transmitter by using DPC. Then, the output SNR at receiver  $\pi(i)$  for data stream  $i$  is given by

$$\zeta_{\pi(i), d_i} = p_i \gamma_{\pi(i), d_i} \quad (8)$$

where  $\gamma_{\pi(i), d_i} = \lambda_{\pi(i), d_i} \beta_i$ , with  $\beta_i = |l_{i,i}|^2$ .

Given the optimal user set  $\pi$  and the corresponding eigen-channel set  $d$ , the sum rate has the form

$$R_{\text{ZFDPC-SUS}} = \max_{p_i: \sum_{i=1}^L p_i \leq P} \sum_{i=1}^L \log_2(1 + p_i \gamma_{\pi(i), d_i}). \quad (9)$$

To maximize (9), the power should be allocated according to the standard water-filling algorithm.

Now consider the problem of selecting the optimal user set  $\pi$  and corresponding eigen-mode index set  $d$ . These sets are chosen to maximize the sum rate, given by (9). When  $M < K$ , to find the optimal solution, one must apply an *exhaustive search* over all possible  $L$ , and for each  $L$ , over all possible sets of  $L$  subchannels taken from the set of  $\sum_{k=1}^K \min\{M, N_k\}$  available eigen-channels spanned by all  $K$  users. Thus, the total number of possible user and eigen-channel selection sets is given by  $\sum_{l=1}^M \binom{\sum_{k=1}^K \min\{M, N_k\}}{l}$ . Further, since different orderings of a given set will yield different output SNRs, all permutations of a given set must also be considered. Clearly, the complexity associated with this exhaustive search is computationally prohibitive in practice, for all but small values of  $K$ .

Here we consider a user and eigen-mode selection algorithm with significantly lower complexity, based on SUS. This algorithm, which was first presented in [13] in the context of ZFBF, iteratively selects a user-eigenmode index pair by searching for a set of users with near orthogonal channel vectors, and is described as follows. Let  $\mathcal{U}_n$  denote the *candidate set* at the  $n$ -th iteration. This set contains the indices of all users and the corresponding eigen-channels that have not been selected previously, and which have not been pruned in the previous iterations (i.e., they have satisfied the ‘‘semi-orthogonality criteria’’ in each of the previous iterations). Also, let  $\mathcal{S}_n = \{(\pi(1), d_1), \dots, (\pi(n), d_n)\}$  denote the set of indices of the selected users and the corresponding eigen-channels after the  $n$ -th iteration.

### ZFDPC-SUS (Algorithm 1)

#### 1) Initialization:

Set  $n = 1$  and  $\mathcal{U}_1 = \{(k, j) \mid k = 1, 2, \dots, K; j = 1, 2, \dots, \min(N_k, M)\}$ .

Let  $\gamma_{k,j}(1) = \lambda_{k,j}$ . The transmitter selects the first user and eigen-channel pair as follows:

$$(\pi(1), d_1) = \arg \max_{(k,j) \in \mathcal{U}_1} \gamma_{k,j}(1). \quad (10)$$

Set  $\mathcal{S}_1 = \{(\pi(1), d_1)\}$ , and define  $\mathbf{q}_1 = \mathbf{v}_{\pi(1), d_1}^H$ .

2) **While**  $n \leq M$ ,  $n \leftarrow n + 1$ .

Calculate candidate set as

$$\begin{aligned} \mathcal{U}_n &= \{(k, j) | (k, j) \in \mathcal{U}_{n-1}, \\ &\quad (k, j) \neq (\pi(n-1), d_{n-1}), |\mathbf{v}_{k,j}^H \mathbf{q}_{n-1}^H|^2 < \delta\} \end{aligned}$$

where  $\delta$  is a positive constant, termed the *semi-orthogonality parameter*, that is preset before the start of the selection procedure.

If  $\mathcal{U}_n$  is empty, set  $n = n - 1$  and go to step 3). Otherwise, for each  $(k, j) \in \mathcal{U}_n$ , denote

$$\xi_i = \mathbf{v}_{k,j}^H \mathbf{q}_i^H, \quad i = 1, \dots, n-1 \quad (11)$$

$$\boldsymbol{\xi}_{k,j} = \mathbf{v}_{k,j}^H - \sum_{i=1}^{n-1} \xi_i \mathbf{q}_i \quad (12)$$

$$\gamma_{k,j}(n) = \lambda_{k,j} \|\boldsymbol{\xi}_{k,j}\|^2. \quad (13)$$

Select the  $n$ -th active user and corresponding eigen-channel as follows:

$$\{(\pi(n), d_n)\} = \arg \max_{(k,j) \in \mathcal{U}_n} \gamma_{k,j}(n). \quad (14)$$

Set

$$\begin{aligned} \mathcal{S}_n &= \mathcal{S}_{n-1} \cup \{(\pi(n), d_n)\}, \\ \mathbf{q}_n &= \frac{\boldsymbol{\xi}_{\pi(n), d_n}}{\|\boldsymbol{\xi}_{\pi(n), d_n}\|}. \end{aligned} \quad (15)$$

3) The transmitter informs the selected users of the indices of their selected eigen-channels; then performs DPC, beamforming, and water-filling power allocation, as discussed previously.

Note that this procedure applies Gram-Schmidt orthogonalization to the ordered rows of  $\boldsymbol{\Xi}_{\pi,d}$ , as described by (11), (12) and (15). As such, it also computes the required transmit precoding matrix in (6).

Observe the following important relations. According to the QR decomposition of  $\boldsymbol{\Xi}_{\pi,d}$ ,

$$\mathbf{v}_{\pi(n), d_n}^H = (\mathbf{v}_{\pi(n), d_n}^H \mathbf{q}_n^H) \mathbf{q}_n + \sum_{j=1}^{n-1} (\mathbf{v}_{\pi(n), d_n}^H \mathbf{q}_j^H) \mathbf{q}_j, \quad (16)$$

and  $l_{n,j} = \mathbf{v}_{\pi(n), d_n}^H \mathbf{q}_j^H$ , for  $j < n$ . With (12),

$$\beta_n = |l_{n,n}|^2 = |\mathbf{v}_{\pi(n), d_n}^H \mathbf{q}_n^H|^2 = \|\boldsymbol{\xi}_{\pi(n), d_n}\|^2. \quad (17)$$

In addition, since  $\|\mathbf{v}_{\pi(n),d_n}\|^2 = 1$  and  $\mathbf{q}_i$ ,  $i = 1, \dots, L$  are orthonormal, it can be easily shown that

$$\sum_{j=1}^n |l_{n,j}|^2 = 1, \quad \text{for } n = 1, 2, \dots, L. \quad (18)$$

### B. Zero-Forcing Beamforming Algorithm

The ZFDPC approach described in the previous section has significantly lower complexity than full (capacity-achieving) DPC, however it is still a nonlinear processing strategy, due to the interference cancelation step. Thus, a common method for reducing complexity even further is to remove the interference cancelation and employ linear processing (linear beamforming). It is well-known, however, that establishing the optimal linear beamforming vectors is a very difficult non-convex optimization problem [19]. Instead, sub-optimal but simple linear processing schemes are usually adopted. Here we will study ZFBF which is one of the most popular linear strategies. Unless otherwise indicated, we will employ the same notational symbols as used in the previous sections.

Let  $\mathbf{C}_{\pi,d}^\dagger$  denote the Moore-Penrose inverse of the equivalent channel matrix  $\mathbf{C}_{\pi,d}$ , i.e.,  $\mathbf{C}_{\pi,d}^\dagger = \mathbf{C}_{\pi,d}^H (\mathbf{C}_{\pi,d} \mathbf{C}_{\pi,d}^H)^{-1}$ , and define  $\tilde{\mathbf{c}}_1, \dots, \tilde{\mathbf{c}}_L$  as the columns of  $\mathbf{C}_{\pi,d}^\dagger$ . For ZFBF, the precoding matrix  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_L]$  is constructed with the beamforming vectors  $\mathbf{w}_i = \frac{\tilde{\mathbf{c}}_i}{\|\tilde{\mathbf{c}}_i\|}$ , for  $i = 1, \dots, L$ . Note that this direct implementation of ZFBF requires the explicit computation of the Moore-Penrose inverse of the channel matrix in order to obtain the beamforming vectors. It has been shown in [18], however, that this direct calculation can be circumvented, thereby significantly reducing the computational complexity. To this end, it is convenient to rewrite the decomposition of  $\mathbf{C}_{\pi,d}$  as  $\mathbf{C}_{\pi,d} = \mathbf{\Lambda}_{\pi,d}^{\frac{1}{2}} \mathbf{L}_{\pi,d} \mathbf{Q}_{\pi,d}$ , where  $\mathbf{\Lambda} = \text{diag}\{\lambda_{\pi(1),d_1}, \dots, \lambda_{\pi(L),d_L}\}$  and  $\mathbf{L}_{\pi,d}$ ,  $\mathbf{Q}_{\pi,d}$  are defined as in Section III-A. Letting  $\mathbf{T}_{\pi,d} = \mathbf{L}_{\pi,d}^{-1} = [\mathbf{t}_1, \dots, \mathbf{t}_L]$ , assuming that  $\mathbf{C}_{\pi,d}$  has full row rank, the Moore-Penrose inverse  $\mathbf{C}_{\pi,d}^\dagger$  can be written as

$$\mathbf{C}_{\pi,d}^\dagger = \mathbf{Q}_{\pi,d}^H \mathbf{L}_{\pi,d}^{-1} \mathbf{\Lambda}_{\pi,d}^{-\frac{1}{2}}. \quad (19)$$

Note that calculating the inverse of  $\mathbf{\Lambda}_{\pi,d}^{\frac{1}{2}}$  is trivial (since it is diagonal), whereas the inverse of  $\mathbf{L}_{\pi,d}$  can be computed using a simple iterative algorithm given in [18, Eq. 11].

For ZFBF, the decoded signal for data stream  $\pi(i)$  is easily shown to be given by

$$\begin{aligned} r_{\pi(i),d_i} &= \sqrt{p_i} \mathbf{c}_{\pi(i),d_i} \mathbf{w}_i x_i + \tilde{n}_{\pi(i),d_i} \\ &= \frac{\sqrt{p_i} \lambda_{\pi(i),d_i}}{\|\mathbf{t}_i\|} x_i + \tilde{n}_{\pi(i),d_i} \end{aligned} \quad (20)$$

with corresponding SNR

$$\varrho_{\pi(i),d_i} = \frac{\lambda_{\pi(i),d_i}}{\|\mathbf{t}_i\|^2}. \quad (21)$$



For the given user set  $\pi$  and the corresponding eigen-channel set  $d$ , the sum rate is given by

$$R_{\text{ZFBF-SUS}} = \max_{p_i: \sum_{i=1}^L p_i \leq P} \sum_{i=1}^L \log_2(1 + p_i \varrho_{\pi(i), d_i}), \quad (22)$$

where the optimal power allocation  $\{p_i\}_{i=1}^L$  is obtained, once again, by applying the waterfilling procedure.

For ZFBF, we consider a user and eigen-channel selection algorithm based on SUS, following the same general procedure as in **Algorithm 1**. Note that SUS has previously been applied to ZFBF in [13]. This algorithm typically assumes that each user is equipped with a single receive antenna, however it extends easily to the multiple receive antenna scenario considered in this paper. One key difference between the algorithms in [11, 13, 18] are the specific methods employed for selecting the “best” user in Step 2 of the algorithm. More specifically, in [13], the same method was applied as in (14), whereas [11] applied a method based on selecting one user at each iteration that results in the largest sum rate when combined with previously selected users. Whilst the latter method can result in larger sum rate, here we will consider the former method for analytically tractability. It has been shown, however, that the difference in sum rate between these two methods is minor [18].

#### IV. SUM RATE ANALYSIS – ASYMPTOTIC $K$

In this section, we investigate the average sum rate of each of the above transceiver structures. For tractability, we make the following assumptions throughout this section:

- (i) For each user, only the principal eigen-channel is considered. As such, we drop the indices for the selected eigen-channels (for example, we use  $\gamma_{\pi(i)}$  instead of  $\gamma_{\pi(i), d_i}$ ).
- (ii) The available power  $P$  is divided equally amongst the active users<sup>3</sup>.

Clearly, the sum rate achieved under these two assumptions will serve as a lower bound to the maximum achievable sum rate. We will also assume that each user has  $N$  antennas, and that there are  $L = M$  data streams.

We will investigate the average sum rate of both scheme discussed in the previous section. We focus on establishing asymptotic results as  $K \rightarrow \infty$ , whilst keeping SNR,  $M$ , and  $N$  fixed.

##### A. ZFDPC-SUS Scheme

To analyze the sum rate of the ZFDPC-SUS system, we require the distribution of the output SNR  $\zeta_{\pi(n)}$ , or alternatively the distribution of  $\gamma_{\pi(n)}$ . Let us first determine the distribution of  $\gamma_k(n)$ ,  $n = 1, \dots, M$ , where  $k$  is an *arbitrary* user selected from the candidate set  $\mathcal{U}_n$ .

<sup>3</sup>Note that in practice the transmit power may be optimized (e.g., according to the water-filling strategy). In such cases, the power allocation depends on the instantaneous channel coefficients and thus changes at the fading rate of the channel, which makes the analysis intractable.

Starting with  $n = 1$ ,  $\gamma_k(1)$ ,  $k = 1, \dots, K$ , are independent and identically distributed (i.i.d.), with

$$\gamma_k(1) = \lambda_{k,\max} \quad (23)$$

where  $\lambda_{k,\max}$  is the maximum eigenvalue of  $\mathbf{H}_k^H \mathbf{H}_k$ , whose probability density function (p.d.f.) and cumulative distribution function (c.d.f.) are known in closed-form and are given as follows [20]:

*Lemma 1:* Let  $\mathbf{H} \sim \mathcal{CN}_{N,M}(\mathbf{0}_{N,M}, \mathbf{I}_N \otimes \mathbf{I}_M)$ . The matrix  $\mathbf{H}^H \mathbf{H}$  is complex Wishart, whose maximum eigenvalue has p.d.f.

$$f_{\max}(x) = \sum_{r=1}^p \sum_{s=q-p}^{(p+q-2r)r} a_{r,s} x^s e^{-rx} \quad (24)$$

and c.d.f.

$$F_{\max}(x) = \sum_{r=1}^p \sum_{s=q-p}^{(p+q-2r)r} \frac{a_{r,s}}{\Gamma^{s+1}} \gamma(s+1, rx) \quad (25)$$

where  $p = \min\{M, N\}$ ,  $q = \max\{M, N\}$ ,  $a_{s,r}$  is a constant (dependent on  $M$  and  $N$ ) which can be computed using the simple numerical method in [21], and  $\gamma(\cdot, \cdot)$  is the lower incomplete gamma function.

For  $n \geq 2$ , evaluating the distribution of  $\gamma_k(n)$ ,  $k \in \mathcal{U}_n$ , is significantly more challenging. Particularly, the “max” operation (10) of Step 1 of the previous iteration (i.e., the  $(n-1)$ -th), and also the semi-orthogonality constraint imposed at Step 2 of the current iteration (i.e., the  $n$ -th) will make the exact distribution of the eigen-channel vectors in  $\mathcal{U}_n$  different from the distributions of the eigen-channel vectors in  $\mathcal{U}_l$ ,  $l \leq n-1$ . More specifically, for  $n \geq 2$ , the eigen-channels for users in the candidate set  $\mathcal{U}_n$  are no longer distributed according to the maximum eigen-channel of a complex Wishart matrix (i.e., for  $k \in \mathcal{U}_n$ ,  $\mathbf{v}_k$  is no longer an isotropically distributed unit vector on the complex unit sphere, and  $\lambda_{k,\max}$  is no longer distributed as the maximum eigenvalue of a complex Wishart matrix).

We see from (13) that  $\gamma_k(n)$  involves the *product* of  $\lambda_{k,\max}$  and the projection variable  $\|\boldsymbol{\xi}_k\|^2$ . For the reasons stated above, the exact distributions of both  $\lambda_{k,\max}$  and  $\|\boldsymbol{\xi}_k\|^2$  for  $k \in \mathcal{U}_n$ ,  $n \geq 2$  are currently unknown and appear very difficult to derive analytically. Fortunately, we can make progress by appealing to the “large-user” regime. In particular, when the number of users in the candidate set  $\mathcal{U}_n$  is large, the problem is greatly simplified by invoking the following key lemma, which shows that removing a finite number of users from  $\mathcal{U}_n$  has negligible impact on the statistical properties of the remaining users. Similar results have also been established previously for different system configurations [11, 13, 18].

*Lemma 2:* At the  $n$ -th iteration,  $2 \leq n \leq M$ , conditioned on the previously selected eigen-channel vectors  $\mathbf{c}_{\pi(1)}, \dots, \mathbf{c}_{\pi(n-1)}$ , the eigen-channel vectors in  $\mathcal{U}_n$  are i.i.d. Furthermore, as the size of the candidate user set  $\mathcal{U}_n$  grows large (i.e.  $\lim_{K \rightarrow \infty} |\mathcal{U}_n| = \infty$ ), conditioned on the previously selected

eigen-channels  $\mathbf{c}_{\pi(1)}, \dots, \mathbf{c}_{\pi(n-1)}$ , the eigen-channel for each user in  $\mathcal{U}_n$  converges in distribution to the distribution of the principal eigen-channel of a complex Wishart matrix.

*Proof:* See Appendix A. ■

Note that our result here differs from that of [18] in both the distribution of the channel vectors and also the user selection algorithm.

Equipped with *Lemma 2*, at the  $n$ -th iteration, from the point of view of the users in  $\mathcal{U}_n$ , the eigen-channel vectors of the selected users in the previous iterations (i.e.,  $\mathbf{c}_{\pi(1)}, \dots, \mathbf{c}_{\pi(n-1)}$ ) appear to be *randomly* selected. Thus, the orthonormal basis  $\mathbf{q}_1, \dots, \mathbf{q}_{n-1}$  (generated from  $\mathbf{c}_{\pi(1)}, \dots, \mathbf{c}_{\pi(n-1)}$ ) appears independent of the eigen-channel vectors of the users in  $\mathcal{U}_n$ . This greatly simplifies the following analysis.

We require the exact distribution of  $\gamma_k(n) = \lambda_{k,\max} \|\boldsymbol{\xi}_k\|^2$ . To this end, the major challenge is to derive the c.d.f. of  $\beta_k(n) = \|\boldsymbol{\xi}_k\|^2$  for an arbitrary user  $k \in \mathcal{U}_n$ , i.e.  $F_{\beta(n)}(x) = \Pr(\beta_k(n) \leq x \mid k \in \mathcal{U}_n)$ . Recalling that  $l_{n,j} = \mathbf{v}_{\pi(n),d_n}^H \mathbf{q}_j^H$  for  $j < n$ , with (17) and (18), we can re-express this c.d.f. as follows:

$$\begin{aligned} F_{\beta(n)}(x) &= \Pr(|\mathbf{v}_k^H \mathbf{q}_n^H|^2 \leq x \mid |\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta, \dots, |\mathbf{v}_k^H \mathbf{q}_{n-1}^H|^2 < \delta) \\ &= \Pr\left(\sum_{i=1}^{n-1} |\mathbf{v}_k^H \mathbf{q}_i^H|^2 \geq 1 - x \mid |\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta, \dots, |\mathbf{v}_k^H \mathbf{q}_{n-1}^H|^2 < \delta\right) \\ &= 1 - \frac{\Pr\left(\sum_{i=1}^{n-1} |\mathbf{v}_k^H \mathbf{q}_i^H|^2 \leq 1 - x, |\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta, \dots, |\mathbf{v}_k^H \mathbf{q}_{n-1}^H|^2 < \delta\right)}{\Pr(|\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta, \dots, |\mathbf{v}_k^H \mathbf{q}_{n-1}^H|^2 < \delta)}. \end{aligned} \quad (26)$$

The denominator,  $\mu_n(\delta) \triangleq \Pr(|\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta, \dots, |\mathbf{v}_k^H \mathbf{q}_{n-1}^H|^2 < \delta)$ , denotes the probability that any arbitrary user  $k \in \{1, \dots, K\}$  will belong to the set  $\mathcal{U}_n$ . Note that this probability has also been considered in the context of ZFBF for the MIMO broadcast channel in [13], where a rather loose lower bound was derived. Here we derive an exact expression which applies for large  $K$ , using an alternative derivation approach. For tractability, our result applies for  $\delta < \frac{1}{M-1}$ , which is easy to establish.

*Lemma 3:* With sufficiently large  $K$  and  $\delta < \frac{1}{M-1}$ , the probability that an arbitrary user  $k \in \{1, \dots, K\}$  belongs to the set  $\mathcal{U}_n$ , for  $n \in \{2, \dots, M\}$ , is given by

$$\begin{aligned} \mu_n(\delta) &= \Pr(|\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta, \dots, |\mathbf{v}_k^H \mathbf{q}_{n-1}^H|^2 < \delta) \\ &= \sum_{k=n-1}^{M-1} \binom{M-1}{k} (-1)^k \left[ \sum_{i=0}^{n-1} \binom{n-1}{i} (-1)^{i,k} \right] \delta^k. \end{aligned} \quad (27)$$

*Proof:* See Appendix B. ■

Note that the term ‘‘sufficiently large’’ in *Lemma 3* implies that  $K$  should be large enough such that:

$$\mathcal{K}_n = |\mathcal{U}_n| \approx K \mu_n(\delta) \quad (28)$$

due to the law of large numbers (LLN). In fact, this also places an additional requirement on  $\delta$ , which

must be selected such that as  $K \rightarrow \infty$ ,  $|\mathcal{U}_n|$  becomes sufficiently large (e.g. such that  $\lim_{K \rightarrow \infty} |\mathcal{U}_n| = \infty$ ). More specifically, since  $\delta < 1$ , by examining (28) and (27) and recalling the condition on  $\delta$  in the lemma statement, we can establish the following design criterion:  $\delta$  should be chosen such that

$$\lim_{K \rightarrow \infty} K\delta^{M-1} = \infty \text{ and } \delta < \frac{1}{M-1}. \quad (29)$$

This implies that any  $\delta$  can be selected, as long as it does not approach zero at a rate of  $1/K^{\frac{1}{M-1}}$  or faster as  $K \rightarrow \infty$ , whilst also meeting the technical condition  $\delta < \frac{1}{M-1}$ . These are very mild conditions which are easy to satisfy (for example, choosing  $\delta$  to be any constant less than  $\frac{1}{M-1}$ ). We further discuss the design implications of selecting  $\delta$  in Section IV-C.

The numerator in (26) can be evaluated using similar methods, which leads to the following result:

*Lemma 4:* Let  $k \in \mathcal{U}_n$ ,  $n \in \{2, \dots, M\}$ , and assume  $\delta$  is chosen to satisfy (29). For sufficiently large  $K$ , the c.d.f. of  $\beta_k(n)$ , given in (26), can be expressed as follows:

$$F_{\beta(n)}(x) = \begin{cases} 0, & x \leq 1 - (n-1)\delta \\ 1 - \frac{\Gamma(M)}{\Gamma(M-n+1)\mu_n(\delta)} \\ \quad \times \int_{t_{n-1}} \cdots \int_{t_1} (1 - \sum_{i=1}^{n-1} t_i)^{M-n} dt_1 \cdots dt_{n-1}, & 1 - (n-1)\delta < x \leq 1 \\ 1, & x > 1 \end{cases} \quad (30)$$

where the integral region is given by  $t_i \in \left[0, \min\{\delta, 1 - x - \sum_{j=i+1}^{n-1} t_j\}\right]$ .

For  $n = 2$ , (30) has the closed-form solution

$$F_{\beta(2)}(x) = \begin{cases} 0 & x \leq 1 - \delta \\ \frac{x^{M-1} - (1-\delta)^{M-1}}{1 - (1-\delta)^{M-1}} & 1 - \delta < x \leq 1 \\ 1 & x > 1 \end{cases} \quad (31)$$

*Proof:* See Appendix C. ■

For arbitrary  $M$  and  $n$ , it is difficult to obtain an exact closed-form solution for this c.d.f. Based on the above lemma, however, we can derive closed-form *upper and lower bounds*, as given by the following:

*Lemma 5:* The c.d.f.  $F_{\beta(n)}(x)$ , for  $n \in \{2, \dots, M\}$ , satisfies  $F_{\bar{\beta}(n)}(x) \leq F_{\beta(n)}(x) \leq F_{\underline{\beta}(n)}(x)$ , with  $F_{\bar{\beta}(n)}(x)$  and  $F_{\underline{\beta}(n)}(x)$  given by (32) and (33)

$$F_{\bar{\beta}(n)}(x) = \begin{cases} 0 & x \leq 1 - (n-1)\delta \\ 1 - \frac{\mu_n(\frac{1-x}{n-1})}{\mu_n(\delta)} & 1 - (n-1)\delta < x \leq 1 \\ 1 & x > 1 \end{cases} \quad (32)$$

and

$$F_{\bar{\beta}(n)}(x) = \begin{cases} 0 & x \leq 1 - (n-1)\delta \\ 1 - \frac{I_{1-x}(n-1, M-n+1)}{\mu_n(\delta)} & 1 - (n-1)\delta < x \leq 1 \\ 1 & x > 1 \end{cases} \quad (33)$$

respectively, where  $\mu_n(\cdot)$  is given by (27) and  $I_x(\cdot, \cdot)$  is the regularized incomplete beta function.

Note that for  $n = 2$ ,  $F_{\beta_k(n)}(x) = F_{\bar{\beta}(n)}(x) = F_{\tilde{\beta}(n)}(x)$ .

*Proof:* See Appendix D. ■

Equipped with *Lemma 5*, and with the help of *Lemma 1*, we may now derive upper and lower bounds on the c.d.f. of  $\gamma_k(n)$ . To establish this result, recall that for an arbitrary user  $k \in \mathcal{U}_n$ ,  $n \geq 2$ , then  $\gamma_k(n) = \lambda_{k, \max} \beta_k(n)$ . Also, define  $\bar{\gamma}_k(n) = \lambda_{k, \max} \bar{\beta}_k(n)$  and  $\tilde{\gamma}_k(n) = \lambda_{k, \max} \tilde{\beta}_k(n)$ , with c.d.f.s  $F_{\bar{\gamma}(n)}(x)$  and  $F_{\tilde{\gamma}(n)}(x)$  respectively.

*Lemma 6:* The c.d.f.  $F_{\gamma(n)}(x)$ , for  $n \in \{2, \dots, M\}$ , satisfies  $F_{\bar{\gamma}(n)}(x) \leq F_{\gamma(n)}(x) \leq F_{\tilde{\gamma}(n)}(x)$ , with  $F_{\bar{\gamma}(n)}(x)$  and  $F_{\tilde{\gamma}(n)}(x)$  given by

$$F_{\bar{\gamma}(n)}(x) = F_{\max}\left(\frac{x}{t}\right) - \frac{1}{\mu_n(\delta)} \sum_{k=n-1}^{M-1} \binom{M-1}{k} (-1)^k \left[ \sum_{i=0}^{n-1} \binom{n-1}{i} (-1)^i \left(\frac{i}{n-1}\right)^k \right] \sum_{r=1}^p \sum_{s=q-p}^{(N+M-2r)r} a_{r,s} \\ \times \sum_{j=0}^k \binom{k}{j} r^{k-j-s-1} (-x)^{k-j} \left[ \Gamma(j-k+s+1, rx) - \Gamma\left(j-k+s+1, \frac{rx}{t}\right) \right]. \quad (34)$$

$$F_{\tilde{\gamma}(n)}(x) = F_{\max}\left(\frac{x}{t}\right) - \frac{1}{\mu_n(\delta)} \sum_{k=0}^{M-n} \binom{M-1}{k} (-1)^k \sum_{r=1}^p \sum_{s=q-p}^{(N+M-2r)r} a_{r,s} \sum_{j=0}^{M-k-1} \binom{M-k-1}{j} r^{M-j-s-2} \\ \times (-x)^{M-j-1} \left[ \Gamma(j+s-M+2, rx) - \Gamma\left(j+s-M+2, \frac{rx}{t}\right) \right]. \quad (35)$$

respectively, where  $F_{\max}(\cdot)$ ,  $p$ ,  $q$  and  $a_{r,s}$  are defined as in *Lemma 1*,  $t = 1 - (n-1)\delta$  and  $\Gamma(\cdot, \cdot)$  denotes the upper incomplete gamma function.

For the case  $n = 2$ ,  $F_{\gamma(n)}(x) = F_{\bar{\gamma}(n)}(x) = F_{\tilde{\gamma}(n)}(x)$ .

*Proof:* See Appendix E. ■

Although not shown due to space limitations, these bounds have been confirmed through simulations.

Recall that our primary aim is to characterize the distribution of  $\zeta_{\pi(n)}$ , or equivalently  $\gamma_{\pi(n)}$  which, from (14), is the maximum of a collection of i.i.d. random variables chosen from  $\mathcal{U}_n$ , with common c.d.f.  $F_{\gamma(n)}(x)$ . Moreover, as discussed previously, our main interest is the case where the number of users  $K$ , and consequently the size of  $\mathcal{U}_n$ , is large. As such, from the theory of extreme order statistics (see e.g. [14, Appendix I] [22]), the asymptotic distribution of the largest order statistic  $\gamma_{\pi(n)}$  depends on the *tail* behavior (large  $x$ ) of  $F_{\gamma(n)}(x)$ . For  $n \geq 2$ , the following closed-form asymptotic (high  $x$ ) expansions for

the c.d.f. upper and lower bounds in (34) and (35) are derived in Appendix F:

$$F_{\tilde{\gamma}(n)}(x) = 1 - \frac{1}{\mu_n(\delta) \epsilon_n} e^{-x} x^{M+N-n-1} + O(e^{-x} x^{M+N-n-2}) \quad (36)$$

$$F_{\bar{\gamma}(n)}(x) = 1 - \frac{1}{\mu_n(\delta) \epsilon_n} e^{-x} x^{M+N-n-1} + O(e^{-x} x^{M+N-n-2}) \quad (37)$$

where

$$\frac{1}{\epsilon_n} = \frac{\Gamma(n)}{\Gamma(M-n+1)\Gamma(N)(n-1)^{n-1}}, \quad (38)$$

$$\frac{1}{\epsilon_n} = \frac{1}{\Gamma(M-n+1)\Gamma(N)}. \quad (39)$$

Based on the above results, we can establish upper and lower bounds of the asymptotic distribution of  $\gamma_{\pi(n)}$ , for large  $K$ . To this end, define  $\tilde{\gamma}_{\pi(n)} = \max_{k \in \mathcal{U}_n} \tilde{\gamma}_k(n)$  and  $\bar{\gamma}_{\pi(n)} = \max_{k \in \mathcal{U}_n} \bar{\gamma}_k(n)$ , with c.d.f.s  $F_{\tilde{\gamma}_{\pi(n)}}(x)$  and  $F_{\bar{\gamma}_{\pi(n)}}(x)$  respectively. It is clear that  $F_{\tilde{\gamma}_{\pi(n)}}(x) \leq F_{\gamma_{\pi(n)}}(x) \leq F_{\bar{\gamma}_{\pi(n)}}(x)$ , where the equalities hold when  $n = 1$ . Then, we have the following lemma:

*Lemma 7:* The random variables  $\tilde{\gamma}_{\pi(n)}$  and  $\bar{\gamma}_{\pi(n)}$ ,  $n \in \{2, \dots, M\}$ , satisfy

$$\begin{aligned} \Pr\{u_n - \log \log \sqrt{K} \leq \tilde{\gamma}_{\pi(n)} \leq u_n + \log \log \sqrt{K}\} \\ \geq 1 - O\left(\frac{1}{\log K}\right), \end{aligned} \quad (40)$$

$$\begin{aligned} \Pr\{\chi_n - \log \log \sqrt{K} \leq \bar{\gamma}_{\pi(n)} \leq \chi_n + \log \log \sqrt{K}\} \\ \geq 1 - O\left(\frac{1}{\log K}\right), \end{aligned} \quad (41)$$

where<sup>4</sup>

$$u_n = \log\left(\frac{K}{\epsilon_n}\right) + (M+N-n-1) \log \log\left(\frac{K}{\epsilon_n}\right), \quad (42)$$

$$\chi_n = \log\left(\frac{K}{\epsilon_n}\right) + (M+N-n-1) \log \log\left(\frac{K}{\epsilon_n}\right). \quad (43)$$

*Proof:* This result is readily established by combining (36) and (37) with the extreme order statistics result given in<sup>5</sup> [14, Lemma 7]. ■

<sup>4</sup>Here  $\log(\cdot)$  represents the natural logarithm.

<sup>5</sup>Note that there are some minor typographical errors with [14, Lemma 7]. Here we have adopted the correct results.

For the case  $n = 1$ ,  $\gamma_{\pi(n)} = \tilde{\gamma}_{\pi(n)} = \bar{\gamma}_{\pi(n)}$ , whose asymptotic distribution is [14]

$$\begin{aligned} \Pr\{u_1 - \log \log \sqrt{K} \leq \gamma_{\pi(1)} \leq u_1 + \log \log \sqrt{K}\} \\ \geq 1 - O\left(\frac{1}{\log K}\right). \end{aligned} \quad (44)$$

Interestingly, we can obtain the same result if we substitute  $n = 1$  into (40)–(43). The asymptotic distribution of  $\zeta_{\pi(n)}$  follows from the above results.

*Lemma 8:* Let  $\rho = \frac{P}{M}$ . For  $\zeta_{\pi(n)}$ ,  $n \in \{1, \dots, M\}$ , we have

$$\begin{aligned} \Pr\{\varpi_n - \rho \log \log \sqrt{K} \leq \zeta_{\pi(n)} \leq v_n + \rho \log \log \sqrt{K}\} \\ \geq 1 - O\left(\frac{1}{\log K}\right), \end{aligned} \quad (45)$$

where

$$\varpi_n = \rho \log \left(\frac{K}{\epsilon_n}\right) + \rho(M + N - n - 1) \log \log \left(\frac{K}{\epsilon_n}\right), \quad (46)$$

$$v_n = \rho \log \left(\frac{K}{\epsilon_n}\right) + \rho(M + N - n - 1) \log \log \left(\frac{K}{\epsilon_n}\right). \quad (47)$$

*Proof:* See Appendix G. ■

We can now prove the following theorem (see Appendix H), which presents a key contribution:

*Theorem 1:* For a fixed number of transmit antennas  $M$  and receive antennas  $N$ , and fixed transmit power  $P$ , if the semi-orthogonality parameter  $\delta$  is chosen to satisfy (29), then the sum rate  $R_{\text{ZFDPC-SUS}}$  of the proposed ZFDPC-SUS scheme satisfies

$$\lim_{K \rightarrow \infty} \frac{R_{\text{ZFDPC-SUS}}}{M \log_2[\rho \log K]} = 1 \quad (48)$$

with probability 1, where  $\rho = P/M$ . In addition,

$$\lim_{K \rightarrow \infty} \mathcal{E}\{R_{\text{BC}}\} - \mathcal{E}\{R_{\text{ZFDPC-SUS}}\} = 0, \quad (49)$$

where  $R_{\text{BC}}$  denotes the sum rate of the MIMO broadcast channel, achieved with DPC. As  $K \rightarrow \infty$ , the average sum rate difference between ZFDPC-SUS and DPC is no greater than  $O\left(\frac{\log \log K}{\log K}\right)$ .

Note that the sum rate difference convergence (49) is much stronger than the sum rate ratio convergence in probability (48), since the latter does not preclude the existence of an infinite sum rate gap between the proposed scheme and the optimal scheme.

### B. ZFBF-SUS Scheme

In this section, we will evaluate the performance of linear ZFBF with SUS. For our analysis, following [13], we will assume that the criterion (14) is used at each iteration of the SUS algorithm to select the best user. In [13], it has been proved that ZFBF-SUS can achieve the same asymptotic sum rate scaling as DPC. Here we establish the stronger result that the average sum rate of ZFBF-SUS converges to the average sum rate achieved with optimal DPC, which was not established in [13]. Deriving an exact expression for the asymptotic distribution of the output SNR for each data stream, analogous to (45), appears very difficult for ZFBF-SUS. Thus, here we adopt a different approach, based on first applying an upper bound which relates the output SNR of ZFBF-SUS in terms of the output SNR of ZFDPC-SUS, and then applying results from the previous subsection. This leads to the following key theorem:

*Theorem 2:* For a fixed number of transmit antennas  $M$  and receive antennas  $N$ , and fixed transmit power  $P$ , if the semi-orthogonality parameter  $\delta$  is chosen to satisfy (29), then the sum rate  $\mathcal{E}\{R_{\text{ZFBF-SUS}}\}$  of the ZFBF-SUS scheme satisfies:

$$\lim_{K \rightarrow \infty} \mathcal{E}\{R_{\text{BC}}\} - \mathcal{E}\{R_{\text{ZFBF-SUS}}\} = 0. \quad (50)$$

As  $K \rightarrow \infty$ , the average sum rate difference between ZFBF-SUS and DPC is no greater than  $O\left(\frac{\log \log K}{\log K}\right)$ .

*Proof:* See Appendix I. ■

This result shows that, as for the ZFDPC-SUS scheme, we can significantly reduce the complexity of the SUS search algorithm by choosing  $\delta$  reasonably small, whilst at the same time achieve the optimal asymptotic sum rate of DPC.

### C. Discussion of Results

Based on the analysis above, some interesting observations are readily in order.

- 1) Asymptotically, both schemes can achieve the maximum spatial multiplexing gain of  $M$ , and also the maximum multi-user diversity gain up to first order (i.e. the SNR scales with  $\log K$ , and the sum rate scales as  $\log \log K$ ). For ZFBF, this scaling behavior agrees with previous results [15, 18].
- 2) As shown in *Theorem 1* and *Theorem 2*, provided that the semi-orthogonality parameter  $\delta$  is selected appropriately, the asymptotic ergodic sum rates of both schemes converge to that of the MIMO broadcast channel, and in both cases the difference in average sum rate with respect to optimal DPC is no greater than  $O\left(\frac{\log \log K}{\log K}\right)$ . Note that similar scaling results have also been obtained for other user selection schemes with ZFBF [15, 18].
- 3) In contrast to most related work, our results provide key insights into the effect of the SUS semi-orthogonality parameter  $\delta$  and the number of receive antennas  $N$ . Considering ZFDPC-SUS, from



(45) and the expressions for  $\varpi_n$  in (46) and  $v_n$  in (47), we see that imposing the constraint  $\delta$  does *not reduce the multi-user diversity gain in both first order terms  $O(\log K)$  and second-order terms  $O(\log \log K)$* . It appears that this result can not be established based on previous (less accurate) SUS analysis methods [13]. Moreover, our analysis demonstrates that whilst the first order terms  $O(\log K)$  in the multi-user diversity gain are unaffected by the number of receive antennas  $N$ , the second-order term grows linearly with both  $N$  and  $M$ . This is consistent with a similar conclusion made in [14], which considered a different system configuration.

- 4) We can also draw insights into the design of  $\delta$ . For practical systems with *finite* numbers of users, obtaining the exact  $\delta$  which yields the optimal complexity–performance tradeoff remains a challenging open problem. However, our asymptotic analysis still provides guidance for the implementation of practical SUS algorithms. In particular, we see that the choice of  $\delta$  is closely related to  $K$  and  $M$  and, to minimize complexity, it is clearly desirable to select  $\delta$  to decrease with increasing  $K$ . At the same time, however, for finite numbers of users it is advisable to “overcompensate” and select  $\delta$  to easily meet the conditions in (29). In our numerical experiments, we found that for systems with  $M \leq 8$ , the choice of  $\delta = \frac{1}{\log K}$  can work well. In addition, since the number of candidate users decreases with each iteration of the SUS algorithm, further complexity savings can be achieved by adaptively selecting  $\delta$ ; e.g., at iteration  $n$ , setting  $\delta_n = \frac{1}{\log |\mathcal{U}_n|}$ .
- 5) Although the results in Section IV-A and IV-B demonstrate that both the ZFDPC-SUS and ZFBF-SUS schemes achieve the same asymptotic average sum rate, the speed of convergence to this optimal sum rate can be very different. Intuitively, this performance difference is caused by a reduction in the *effective channel gain* [13] seen by the ZFBF receivers. Thus, for finite  $K$ , there will be a gap in the average sum rates of the two schemes. We will now study this more closely.

## V. SUM RATE ANALYSIS – FINITE $K$

In this section, we analyze the achievable sum rates of the ZFDPC-SUS and ZFBF-SUS schemes for *finite* numbers of users. To obtain clear insights, we focus on the high and low SNR regimes. Our analysis is based on studying the gap between the sum rates achieved by the two transceivers and a fixed upper bound. This study follows the method of [23], which considered single-user MIMO receivers. We will first evaluate the performance for a given set of channel realizations, and then investigate the average performance via simulations. We make the same assumptions as stated at the beginning of Section IV.

Given a set of  $M$  users  $\pi$  determined by user selection<sup>6</sup>, the sum capacity of the MIMO broadcast chan-

<sup>6</sup>For a meaningful comparison, we will assume that for both schemes, the same SUS selection criteria is used, based on (14). As such, the active users sets and the corresponding compound channel matrix  $\mathbf{C}_{\pi,d}$  will be the same for both schemes.

nel  $\{\mathbf{H}_{\pi(k)}\}_{k=1}^M$  can be written by using the duality of the MIMO broadcast channel and the MIMO multiple access channel as [4]  $C_{\text{BC}}(\{\mathbf{H}_{\pi(k)}\}_{k=1}^M, P) = \max_{\sum_k \text{tr} \mathbf{Q}_k \leq P} \log_2 \det \left( \mathbf{I} + \sum_{k=1}^M \mathbf{H}_{\pi(k)}^H \mathbf{Q}_k \mathbf{H}_{\pi(k)} \right)$ . Since no closed-form solution exists, it is very difficult to compare  $C_{\text{BC}}(\{\mathbf{H}_{\pi(k)}\}_{k=1}^M, P)$  with  $R_{\text{ZFDPC-SUS}}$  and  $R_{\text{ZFBF-SUS}}$ . In fact, even with our assumption of equal power allocation, i.e.  $\mathbf{Q}_k = \frac{P}{KN} \mathbf{I}$ , this problem is still difficult, due to the complicated structure of the compound channel matrix  $\mathbf{C}_{\pi,d}$  for the ZFDPC and ZFBF schemes (see (5)). Thus, to analyze the difference in sum rate between  $R_{\text{ZFDPC-SUS}}$  and  $R_{\text{ZFBF-SUS}}$  for finite  $K$ , we adopt an indirect approach and focus on characterizing the differences between the sum rates achieved by the two transceiver structures and  $C$ , where  $C = \log_2 \det(\mathbf{I}_M + \rho \mathbf{C}_{\pi,d} \mathbf{C}_{\pi,d}^H)$  with  $\rho = P/M$ .

Before presenting our main results, it is worth noting that [5, *Theorem 3*]  $\lim_{P \rightarrow \infty} C_{\text{BC}}(\mathbf{C}_{\pi,d}, P) - C = 0$ , where  $C_{\text{BC}}(\mathbf{C}_{\pi,d}, P)$  denotes the sum capacity of a MIMO broadcast system given by (5). Moreover, for the case  $N = 1$ ,  $\{\mathbf{H}_{\pi(k)}\}_{k=1}^M$  reduces to  $\mathbf{C}_{\pi,d}$  and  $C_{\text{BC}}(\{\mathbf{H}_{\pi(k)}\}_{k=1}^M, P)$  coincides with  $C_{\text{BC}}(\mathbf{C}_{\pi,d}, P)$ . Thus, the high SNR results which we establish below correspond precisely to the gaps between the sum rates achieved by the two transceivers and the sum capacity achieved with optimal DPC. Define

$$\eta_i = \sum_{j=1}^{i-1} \frac{|l_{i,j}|^2}{|l_{i,i}|^2}, \quad \kappa_i = \sum_{j=i+1}^M \frac{|t_{j,i}|^2}{|t_{i,i}|^2}, \quad (51)$$

where  $l_{i,j}$  and  $t_{i,j}$  are the  $(i, j)$ -th elements of matrices  $\mathbf{L}_{\pi,d}$  and  $\mathbf{T}_{\pi,d}$ , respectively. Some basic manipulations of the results in [23] yield the following theorem:

*Theorem 3:* For finite number of users  $K$ , finite number of transmit and receive antennas  $M$  and  $N$ ,

- In the high SNR region:

$$\begin{aligned} C - R_{\text{ZFDPC-SUS}} &= \frac{1}{\rho \log 2} \sum_{i=1}^M \frac{\kappa_i}{\lambda_{\pi(i)} |l_{i,i}|^2} \\ &\quad + O(\rho^{-2}), \end{aligned} \quad (52)$$

$$\begin{aligned} C - R_{\text{ZFBF-SUS}} &= \sum_{i=1}^M \log_2(1 + \kappa_i) \\ &\quad + O(\rho^{-2}). \end{aligned} \quad (53)$$

- In the low SNR region:

$$\begin{aligned} C - R_{\text{ZFDPC-SUS}} &= \frac{\rho}{\log 2} \sum_{i=1}^M \eta_i \lambda_{\pi(i)} |l_{i,i}|^2 \\ &\quad + O(\rho^2), \end{aligned} \quad (54)$$

$$\begin{aligned}
C - R_{\text{ZFBF-SUS}} &= \frac{\rho}{\log 2} \sum_{i=1}^M \left(1 + \eta_i - \frac{1}{1 + \kappa_i}\right) \\
&\quad \times \lambda_{\pi(i)} |l_{i,i}|^2 + O(\rho^2).
\end{aligned} \tag{55}$$

From these results, we can make the following conclusions.

*High SNR Region:* As  $\rho \rightarrow \infty$ , for ZFDPC-SUS the sum rate approaches  $C$ , whereas for ZFBF-SUS there is a constant sum rate gap of  $\mathcal{A} \triangleq \sum_{i=1}^M \log_2(1 + \kappa_i)$ . This gap can be zero only when  $\kappa_i = 0$ , which is a rare case corresponding to complete orthogonality between the row vectors of  $\mathbf{C}_{\pi,d}$ . Subtracting (54) from (55), in this region we can also quantify the sum rate gap between ZFDPC-SUS and ZFBF-SUS as  $R_{\text{ZFDPC-SUS}} - R_{\text{ZFBF-SUS}} = \mathcal{A} + O(\rho^{-1})$ , which shows the advantage of ZFDPC-SUS for finite  $K$ .

*Low SNR Region:* As  $\rho \rightarrow 0$ , for both ZFDPC-SUS and ZFBF-SUS, the sum rate gaps w.r.t.  $C$  approach zero linearly with  $\rho$ . Moreover, in this region we can again quantify the sum rate gap as  $R_{\text{ZFDPC-SUS}} - R_{\text{ZFBF-SUS}} = \frac{\rho}{\log 2} \sum_{i=1}^M \left(1 - \frac{1}{1 + \kappa_i}\right) \lambda_{\pi(i)} |l_{i,i}|^2$ , which is non-negative. It is also worth noting that in the low SNR regime, better performance may be achievable by transmitting with full power to only a single user, rather than sending equal power streams to  $M$  selected users. The benefit of this approach, however, will depend not only on the SNR value, but also on the number of users  $K$ . In particular, the benefit of this approach is expected to be most evident when  $K$  is small, for which case there will be the most disparity between the dominant eigen-channels of the users.

*Effect of SUS Parameter  $\delta$ :* According to the SUS algorithm, we have  $|l_{i,j}|^2 < \delta$  for  $i > j$ , and  $|l_{i,i}|^2 > 1 - (i - 1)\delta$ . Thus, with smaller semi-orthogonality parameter  $\delta$ , it is more likely to have off-diagonal elements with smaller absolute value in both  $\mathbf{L}_{\pi,d}$  and  $\mathbf{T}_{\pi,d}$  (i.e smaller  $|l_{i,j}|, i < j$  and  $|t_{j,i}|, i < j$ ) and more likely to have diagonal elements with larger absolute value in  $\mathbf{L}_{\pi,d}$ . From (51), these observations imply that a smaller  $\delta$  leads to smaller  $\eta_i$  and  $\kappa_i$ . In addition, it is easy to see that  $\eta_i |l_{i,i}|^2 = \sum_{j=1}^{i-1} |l_{i,j}|^2$  and  $(1 + \eta_i) |l_{i,i}|^2 = 1$ . With these results, we see that by decreasing  $\delta$ , the sum rate gaps for both transceivers are likely to decrease, for both high and low SNRs. This implies that the sum rates of both transceivers are likely to increase, which agrees with intuition.

Fig. 1 demonstrates the average sum rate gaps of ZFDPC-SUS and ZFBF-SUS for different SNRs. Results are shown for  $M = 4$ ,  $N = 4$ ,  $K = 50$ , and  $\delta = \frac{1}{\log K}$ . These results confirm our analytical conclusions given above, based on *Theorem 3*.

## VI. NUMERICAL RESULTS

For our simulations, we use  $P = 15$  dB,  $\delta = \frac{1}{\log K}$ , and the optimal water-filling power allocation.

Fig. 2 plots the average sum rate achieved by ZFDPC-SUS and ZFBF-SUS as a function of the number of users. Curves are also presented for ZFBF with complete search, as well as optimal DPC. In the first

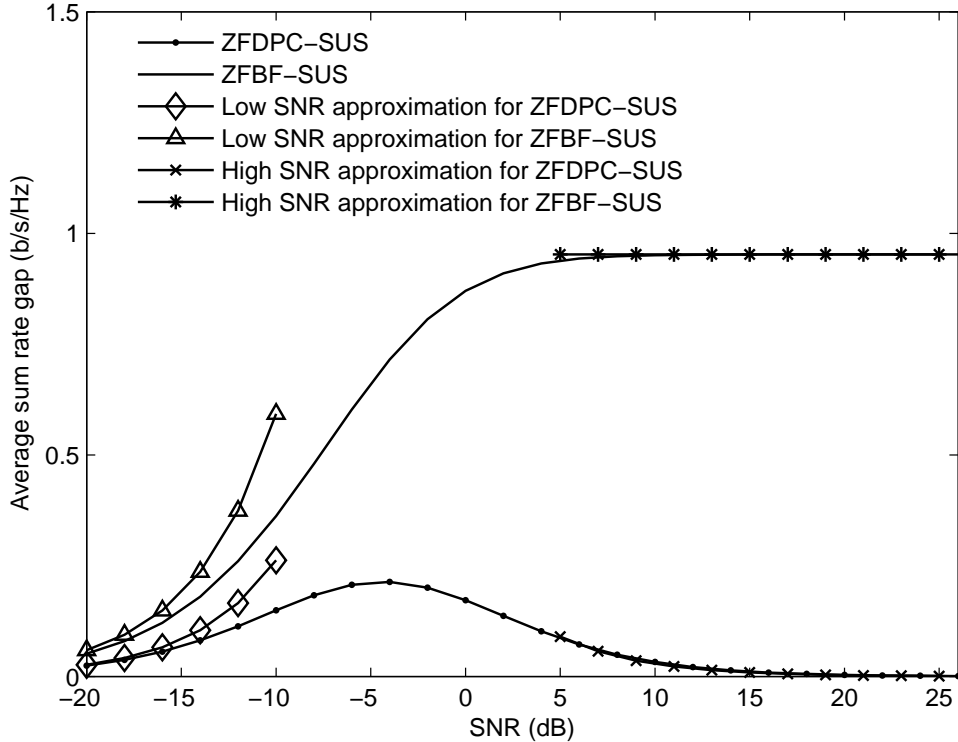


Fig. 1. Comparison of sum rate gap for different SNRs.  $M = 4$ ,  $N = 4$ ,  $K = 50$ .

case, a search is conducted over all combinations of users, and the combination with the highest sum rate is selected. Due to the very high complexity of this approach, we only provide results for relatively small  $K$ . The optimal DPC curve acts as an achievable upper bound, and is computed using the algorithm from [24]. In addition, based on (98) and the expressions for  $u_n$  in (42) and  $\chi_n$  in (43), we have plotted  $\sum_{i=1}^M \log_2(1 + \rho(\log K + (M + N - i - 1) \log \log K))$  as an asymptotic approximation for the average sum rate of the ZFDPC-SUS scheme. As evident from the figure, the performance of ZFDPC-SUS is very close to that of DPC, and is slowly converging to DPC as  $K$  grows large. The asymptotic approximation for ZFDPC-SUS based on our analysis is also quite good (within 1 bps/Hz). Considering ZFBF, we see that the ZFBF-SUS curve is no more than 0.5 dB away from that of the complete search method; further verifying the utility of the SUS approach. Moreover, the ZFBF curves are far below the ZFDPC-SUS curve, demonstrating that ZFDPC-SUS has *significant* performance advantages at finite  $K$ . For further comparison, we have also implemented a related algorithm proposed in [15] and plotted the corresponding sum rate curve. This curve is generated by using an optimal threshold, computed by an exhaustive search. The performance is close to that of ZFBF-SUS.

Fig. 3 compares the average sum rate of ZFDPC-SUS and ZFBF-SUS as a function of the number of

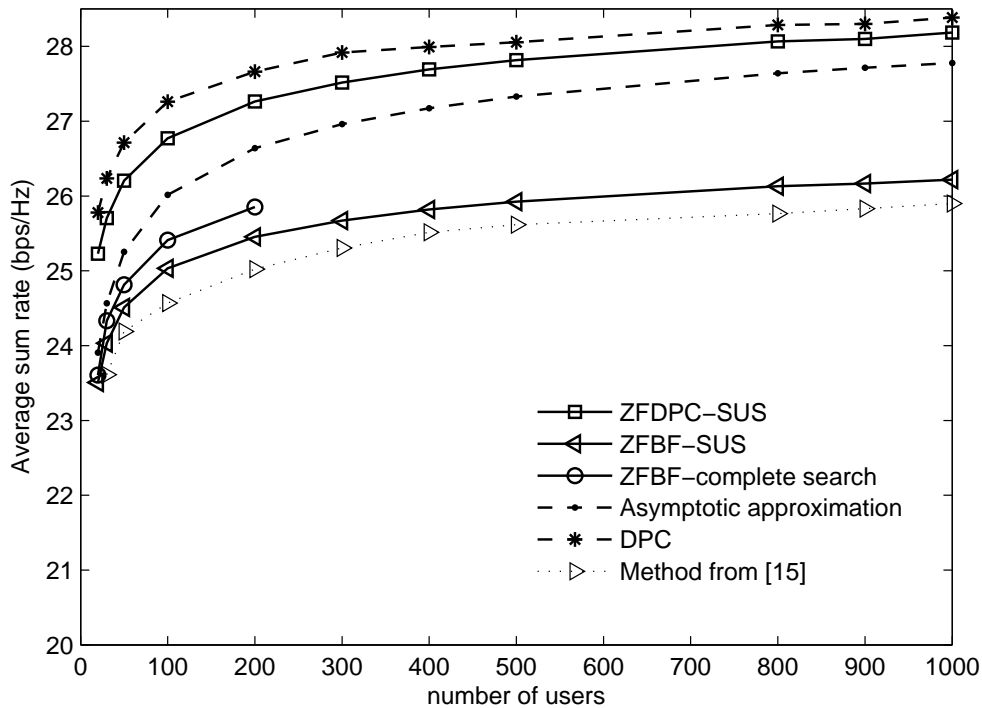


Fig. 2. Comparison of average sum rates for different numbers of users.  $M = 4$ ,  $N = 4$ ,  $P = 15$  dB.

users, for different numbers of receive antennas. Note that according to (98) and the expressions for  $u_n$  and  $\chi_n$  in (42) and (43) respectively, if we increase the number of receive antennas by one, the increase in sum rate can be approximated as  $M \log \left( 1 + \frac{\rho \log \log K}{1 + \rho \log K} \right) \rightarrow 0$  as  $K \rightarrow \infty$ ; i.e., the difference in sum rate will be negligible for large  $K$ . However, the figure shows that this convergence is very slow, and that increasing the number of receive antennas can significantly increase the sum rate for finite  $K$ .

## VII. CONCLUSION

We have investigated the sum rate of two low complexity eigenmode-based transmission techniques for the MIMO broadcast channel, ZFDPC-SUS and ZFBF-SUS. We proved that ZFDPC-SUS can achieve the optimal sum rate scaling of the MIMO broadcast channel, and that the average sum rate of both techniques converges to the average sum capacity of the MIMO broadcast channel as  $K$  grows large (albeit at different rates). We also investigated and compared the achievable sum rates of ZFDPC-SUS and ZFBF-SUS for finite  $K$ , and demonstrated that ZFDPC-SUS has significant performance advantages. In contrast to most previous related results, our analytical results provide important insights into the benefit of multiple receive antennas, and the effect of the SUS algorithm.

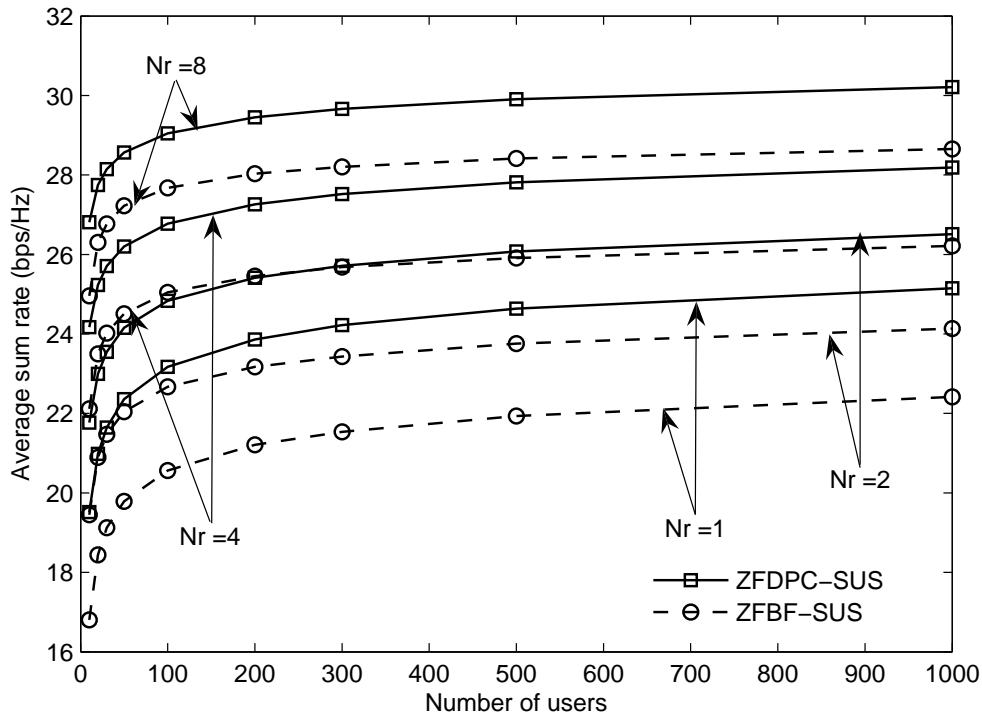


Fig. 3. Comparison of average sum rates for different numbers of users and different numbers of receive antennas.  $M = 4$ ,  $P = 15$  dB.

## APPENDIX A

### PROOF OF Lemma 2

Our derivation closely follows the method of proof for [18, Lemma 3] and [25, Lemma 1]. For two complex vectors  $\mathbf{z} = \mathbf{z}_r + j\mathbf{z}_i$  and  $\mathbf{z}' = \mathbf{z}'_r + j\mathbf{z}'_i$  with the same dimension, we write  $\mathbf{z} \preceq \mathbf{z}'$  if every element of  $\mathbf{z}_r$  and  $\mathbf{z}_i$  is less than or equal to its counterpart in  $\mathbf{z}'_r$  and  $\mathbf{z}'_i$ , respectively. Let  $\mathcal{K}_n$  denote the cardinality of the candidate set  $\mathcal{U}_n$ . For the first iteration,  $\mathcal{K}_1 = K$  and  $\mathbf{c}_{\pi(1)}$  is the vector with the maximum norm. For clarity of exposition, at the end of  $n$ -th iteration, we relabel the eigen-channel vectors in  $\mathcal{U}_n/\{\pi(n)\}$  as  $\tilde{\mathbf{c}}_1, \dots, \tilde{\mathbf{c}}_{\mathcal{K}_n-1}$ .

We find that the result in [25, Lemma 1], which was derived specifically for Gaussian vectors, holds more generally and does not require the Gaussian assumption, and indeed can also be adapted to our case. The proof is based on induction. For the first iteration, we have

$$\begin{aligned} & \Pr\{\tilde{\mathbf{c}}_1 \preceq \mathbf{z}_1, \dots, \tilde{\mathbf{c}}_{K-1} \preceq \mathbf{z}_{K-1} | \mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}\} \\ &= \prod_{i=1}^{K-1} \Pr\{\tilde{\mathbf{c}}_i \preceq \mathbf{z}_i | \|\tilde{\mathbf{c}}_i\| < \|\mathbf{z}_{(1)}\|\} \end{aligned} \quad (56)$$

and since  $\lim_{K \rightarrow \infty} \|\mathbf{z}_{(1)}\| = \infty$ ,

$$\lim_{K \rightarrow \infty} \Pr\{\tilde{\mathbf{c}}_i \preceq \mathbf{z}_i \mid \|\tilde{\mathbf{c}}_i\| < \|\mathbf{z}_{(1)}\|\} = F_{\mathbf{c}}(\mathbf{z}_i), \quad (57)$$

where  $F_{\mathbf{c}}(\cdot)$  is the c.d.f. of the principal eigen-vector of a complex Wishart matrix.

Now assume that this lemma holds up to the  $(n-1)$ -th iteration and let us consider the  $n$ -th iteration. Conditioned on  $\mathbf{c}_{\pi(1)}, \dots, \mathbf{c}_{\pi(n-1)}$ , according to our assumption, the channel vectors in  $\mathcal{U}_n$  are i.i.d. and converge in distribution to the principal eigen-vector of a complex Wishart matrix. At the end of step 3) of the  $n$ -th iteration, user  $\pi(n)$  is chosen. Any user  $k$  in  $\mathcal{U}_n$  satisfies  $\gamma_k(n) \leq \gamma_{\pi(n)}$ . Replacing the condition<sup>7</sup>  $\{\mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}\}$  and  $\{\|\tilde{\mathbf{c}}_i\| \leq \|\mathbf{z}_{(1)}\|\}$  by  $\{\mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{c}_{\pi(n-1)} = \mathbf{z}_{(n-1)}, \mathbf{c}_{\pi(n)} = \mathbf{z}_{(n)}\}$  and  $\{\mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{c}_{\pi(n-1)} = \mathbf{z}_{(n-1)}, \gamma_k(n) \leq \gamma_{\pi(n)}\}$  respectively in the derivation in [25, Lemma 1] and following the same method as in [25, Lemma 1], we can establish that the remaining channel vectors in  $\mathcal{U}_n$  are i.i.d. with c.d.f.

$$\Pr\{\tilde{\mathbf{c}}_i \preceq \mathbf{z}_i \mid \mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{c}_{\pi(n-1)} = \mathbf{z}_{(n-1)}, \gamma_k(n) \leq \gamma_{\pi(n)}\} \quad (58)$$

for  $i = 1, \dots, \mathcal{K}_n - 1$ . Since  $\lim_{K \rightarrow \infty} \mathcal{K}_n = \infty$ ,  $\gamma_{\pi(n)}$  is unbounded from above, i.e.,

$$\lim_{K \rightarrow \infty} \gamma_{\pi(n)} = \infty, \quad (59)$$

and we have

$$\begin{aligned} & \lim_{K \rightarrow \infty} \Pr\{\tilde{\mathbf{c}}_i \preceq \mathbf{z}_i \mid \mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{c}_{\pi(n-1)} = \mathbf{z}_{(n-1)}, \gamma_k(n) \leq \gamma_{\pi(n)}\} \\ &= \Pr\{\tilde{\mathbf{c}}_i \preceq \mathbf{z}_i \mid \mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{c}_{\pi(n-1)} = \mathbf{z}_{(n-1)}\}. \end{aligned} \quad (60)$$

By induction  $\Pr\{\tilde{\mathbf{c}}_i \preceq \mathbf{z}_i \mid \mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{c}_{\pi(n-1)} = \mathbf{z}_{(n-1)}\}$  converges in distribution to the distribution of the principal eigen-vector of a complex Wishart matrix, thereby establishing the lemma.

<sup>7</sup>To be more precise, we note that different notation is used in [18]. Our conditions  $\{\mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{c}_{\pi(n)} = \mathbf{z}_{(n)}\}$  and  $\{\mathbf{c}_{\pi(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{c}_{\pi(n-1)} = \mathbf{z}_{(n-1)}, \gamma_k(n) \leq \gamma_{\pi(n)}\}$  are analogous to the conditions  $\{\mathbf{h}_{j(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{h}_{j(n)} = \mathbf{z}_{(n)}\}$  and  $\{\mathbf{h}_{j(1)} = \mathbf{z}_{(1)}, \dots, \mathbf{h}_{j(n-1)} = \mathbf{z}_{(n-1)}, R_{(n)}^{\text{BF}}(\mathbf{h}_i) \leq R_{(n)}^{\text{BF}}(\mathbf{z}_{(n)})\}$  given in [18].

## APPENDIX B

## PROOF OF Lemma 3

According to Lemma 2, the eigen-vector  $\mathbf{v}_k$ , for  $k \in \mathcal{U}_n$ , is an isotropically distributed unit vector on the  $M$ -dimensional complex unit hypersphere. In addition, for large  $K$ , the subspace spanned by the orthonormal basis  $\mathbf{q}_1, \dots, \mathbf{q}_{n-1}$  becomes independent of  $\mathbf{v}_k$ . Thus, without loss of generality we can assume  $\mathbf{q}_i = \mathbf{e}_i$ , where  $\mathbf{e}_i$  is the  $i$ -th row of the identity matrix  $\mathbf{I}_M$ . Let  $\mathbf{v}_k = [v_1, \dots, v_M]^T$ , then

$$\begin{aligned} \mu_n(\delta) &= \Pr(|\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta, \dots, |\mathbf{v}_k^H \mathbf{q}_{n-1}^H|^2 < \delta) \\ &= \Pr(|v_1|^2 < \delta, \dots, |v_{n-1}|^2 < \delta). \end{aligned} \quad (61)$$

In the following we will first derive the joint p.d.f. of  $|v_1|^2, \dots, |v_{n-1}|^2$ .

The surface area of a complex unit hypersphere of  $M$  dimensions is  $\frac{2\pi^M}{\Gamma(M)}$  [26]. So the joint p.d.f. of  $v_1, \dots, v_M$  can be written as:

$$f(\mathbf{v}_k) = f(v_1, \dots, v_M) = \begin{cases} \frac{\Gamma(M)}{2\pi^M}, & \|\mathbf{v}_k\| = 1 \\ 0, & \text{otherwise} \end{cases}. \quad (62)$$

Define  $v_i = x_{2i-1} + jx_{2i}$ . Then, the joint p.d.f. of  $x_1, \dots, x_{2M}$  can be expressed as:

$$f(x_1, x_2, \dots, x_{2M}) = \begin{cases} \frac{\Gamma(M)}{2\pi^M}, & \sum_{i=1}^{2M} x_i^2 = 1 \\ 0, & \text{otherwise} \end{cases}. \quad (63)$$

We require the joint p.d.f. of  $x_1, \dots, x_{2(n-1)}$ , which is evaluated via

$$\begin{aligned} &f(x_1, \dots, x_{2(n-1)}) \\ &= \int \cdots \int_{\sum_{i=1}^{2M} x_i^2 = 1} f(x_1, \dots, x_{2M}) \\ &\quad \times dx_{2(n-1)+1} \cdots dx_{2M} \\ &= \frac{\Gamma(M)}{2\pi^M} V(x_1, \dots, x_{2(n-1)}) \end{aligned} \quad (64)$$

where  $V(x_1, \dots, x_{2(n-1)})$  denotes the area

$$\begin{aligned} &V(x_1, \dots, x_{2(n-1)}) \\ &= \int \cdots \int_{\sum_{i=1}^{2M} x_i^2 = 1} dx_{2(n-1)+1} \cdots dx_{2M} \\ &= \int \cdots \int_{\sum_{i=2(n-1)+1}^{2M} x_i^2 = 1 - \sum_{i=1}^{2(n-1)} x_i^2} \\ &\quad \times dx_{2(n-1)+1} \cdots dx_{2M}. \end{aligned} \quad (65)$$



The multi-dimensional integral (65) is seen to be the surface area of a real  $(2M - 2(n - 1))$ -dimensional hypersphere of radius  $\sqrt{1 - \sum_{i=1}^{2(n-1)} x_i^2}$ . Thus, using results from [26], we evaluate this integral as follows:

$$\begin{aligned} V(x_1, \dots, x_{2(n-1)}) &= \frac{2\pi^{M-n+1}}{\Gamma(M-n+1)} \left(1 - \sum_{i=1}^{2(n-1)} x_i^2\right)^{\frac{2(M-n+1)-1}{2}} \\ &\quad \times \sqrt{\det \mathbf{A}} \, dx_1 \cdots dx_{2(n-1)}, \end{aligned} \quad (66)$$

where  $\mathbf{A}$  is a  $(2(n-1) + 1) \times (2(n-1) + 1)$  matrix with  $(i, j)$ -th element  $\mathbf{A}_{i,j} = \frac{\partial \boldsymbol{\theta}}{\partial x_i} \cdot \frac{\partial \boldsymbol{\theta}}{\partial x_j}$  with  $\boldsymbol{\theta} = \left(x_1, \dots, x_{2(n-1)}, \sqrt{1 - \sum_{i=1}^{2(n-1)} x_i^2}\right)^T$ , and ‘ $\cdot$ ’ denotes the vector inner product operation. We can compute  $\mathbf{A}_{i,j} = \delta_{i,j} + \frac{x_i x_j}{1 - \sum_{i=1}^{2(n-1)} x_i^2}$ , where  $\delta_{i,j}$  is the Kronecker-delta function, and after some manipulations obtain  $\det \mathbf{A} = \frac{1}{1 - \sum_{i=1}^{2(n-1)} x_i^2}$ . Combining this result with (64) and (66) we obtain

$$\begin{aligned} f(x_1, \dots, x_{2(n-1)}) &= \frac{\Gamma(M)}{\Gamma(M-n+1)\pi^{n-1}} \\ &\quad \times \left(1 - \sum_{i=1}^{2(n-1)} x_i^2\right)^{M-n}. \end{aligned} \quad (67)$$

It is now convenient to make the polar coordinate transformations  $x_{2i-1} = r_i \cos \theta_i$ ,  $x_{2i} = r_i \sin \theta_i$ , for  $i = 1, \dots, n-1$ , where  $r_i \geq 0$ ,  $0 \leq \theta_i \leq 2\pi$ . The corresponding Jacobian is easily evaluated as [26]  $\left(\prod_{i=1}^{n-1} r_i\right)^{-1}$ . So the joint density of  $r_1, \dots, r_{n-1}$  is

$$\begin{aligned} f(r_1, \dots, r_{n-1}) &= \frac{\Gamma(M)}{\Gamma(M-n+1)\pi^{n-1}} \left(1 - \sum_{i=1}^{n-1} r_i^2\right)^{M-n} \prod_{i=1}^{n-1} r_i \\ &\quad \times \prod_{i=1}^{n-1} \int_0^{2\pi} d\theta_i \\ &= \frac{2^{n-1}\Gamma(M)}{\Gamma(M-n+1)} \left(1 - \sum_{i=1}^{n-1} r_i^2\right)^{M-n} \prod_{i=1}^{n-1} r_i. \end{aligned} \quad (68)$$

Next we apply the transformation  $t_i = r_i^2$ ,  $i = 1, \dots, n-1$ . Clearly  $t_i = |v_i|^2$  (we will deal with  $t_i$  subsequently to simplify notation). The corresponding Jacobian is  $J(t_1, \dots, t_{n-1}) = 1/(2^{n-1}\sqrt{t_1, \dots, t_{n-1}})$ . So we obtain the desired joint p.d.f. of  $t_1, \dots, t_{n-1}$  as

$$f(t_1, \dots, t_{n-1}) = \frac{\Gamma(M)}{\Gamma(M-n+1)} \left(1 - \sum_{i=1}^{n-1} t_i\right)^{M-n}. \quad (69)$$

Armed with this result, we can now evaluate the desired probability  $\mu_n(\delta)$  in (61). For notational

convenience, we will consider  $\mu_{n+1}(\delta)$ , for  $n+1 \in \{2, \dots, M\}$ . Denoting  $D_n = \{0 \leq t_1 \leq \delta, \dots, 0 \leq t_n \leq \delta\}$ , we have

$$\begin{aligned}\mu_{n+1}(\delta) &= \int \cdots \int_{D_n} f(t_1, \dots, t_n) dt_1 \cdots dt_n \\ &= \frac{\Gamma(M)}{\Gamma(M-n)} \varphi_n(1)\end{aligned}\quad (70)$$

where we have defined

$$\varphi_n(z) = \int \cdots \int_{D_n} \left( z - \sum_{i=1}^n t_i \right)^{M-n-1} dt_1 \cdots dt_n \quad (71)$$

for  $z \geq n\delta$ . Note that with this definition,  $\varphi_n(1)$  exists for all  $n$  provided that  $\delta < \frac{1}{M-1}$ . This condition is assumed in the lemma statement. Then  $\varphi_n(z)$  can be written as

$$\begin{aligned}\varphi_n(z) &= \int \cdots \int_{D_{n-1}} \left( \int_0^\delta \left( z - \sum_{i=1}^n t_i \right)^{M-n-1} dt_n \right) dt_1 \cdots dt_{n-1} \\ &= \frac{1}{M-n} \int \cdots \int_{D_{n-1}} \left[ \left( z - \sum_{i=1}^{n-1} t_i \right)^{M-n} - \left( z - \delta - \sum_{i=1}^{n-1} t_i \right)^{M-n} \right] dt_1 \cdots dt_{n-1} \\ &= \frac{1}{M-n} (\varphi_{n-1}(z) - \varphi_{n-1}(z - \delta)).\end{aligned}\quad (72)$$

So we have

$$\varphi_n(1) = \frac{1}{M-n} (\varphi_{n-1}(1) - \varphi_{n-1}(1 - \delta)) \quad (73)$$

$$\begin{aligned}&= \frac{1}{(M-n)(M-n+1)} \\ &\times (\varphi_{n-2}(1) - 2\varphi_{n-2}(1 - \delta) + \varphi_{n-2}(1 - 2\delta)).\end{aligned}\quad (74)$$

We will now prove, using mathematical induction, that for any integer  $k \in \{1, 2, \dots, n-1\}$ ,

$$\begin{aligned}\varphi_n(1) &= \left[ \prod_{j=0}^{k-1} (M-n+j) \right]^{-1} \\ &\times \sum_{i=0}^k (-1)^i \binom{k}{i} \varphi_{n-k}(1 - i\delta).\end{aligned}\quad (75)$$

According to (73) and (74), (75) holds for  $k=1$  and  $k=2$  respectively. Assuming that (75) holds for

integer  $k$ , applying (72) in (75) yields

$$\varphi_n(1) = \left[ \prod_{j=0}^k (M - n + j) \right]^{-1} \sum_{i=0}^k (-1)^i \binom{k}{i} \left[ \varphi_{n-k-1}(1 - i\delta) - \varphi_{n-k-1}(1 - (i+1)\delta) \right] \quad (76)$$

$$\begin{aligned} &= \left[ \prod_{j=0}^k (M - n + j) \right]^{-1} \left\{ \varphi_{n-k-1}(1) + (-1)^{k+1} \varphi_{n-k-1}(1 - (k+1)\delta) \right. \\ &\quad \left. + \sum_{i=0}^{k-1} (-1)^{i+1} \binom{k+1}{i+1} \varphi_{n-k-1}(1 - (i+1)\delta) \right\} \end{aligned} \quad (77)$$

$$= \left[ \prod_{j=0}^k (M - n + j) \right]^{-1} \sum_{i=0}^{k+1} (-1)^i \binom{k+1}{i} \varphi_{n-k-1}(1 - i\delta) \quad (78)$$

where, to obtain (77), we have used  $\binom{k}{i+1} = \binom{k-1}{i} + \binom{k-1}{i+1}$ . Thus, from (78), if (75) holds for integer  $k$ , it also holds for  $k+1$ . By induction, (75) then holds for any integer  $1 \leq k < n$ . Setting  $k = n-1$  in (75),

$$\begin{aligned} \varphi_n(1) &= \left[ \prod_{j=0}^{n-2} (M - n + j) \right]^{-1} \\ &\quad \times \sum_{i=0}^{n-1} (-1)^i \binom{n-1}{i} \varphi_1(1 - i\delta). \end{aligned} \quad (79)$$

The function  $\varphi_1(1 - i\delta)$  can be evaluated as

$$\begin{aligned} \varphi_1(1 - i\delta) &= \int_0^\delta (1 - i\delta - t_1)^{M-2} dt_1 \\ &= \frac{(1 - i\delta)^{M-1} - (1 - (i+1)\delta)^{M-1}}{M-1}. \end{aligned} \quad (80)$$

Substituting (80) into (79) yields a closed-form solution, which we simplify as follows:

$$\begin{aligned} \varphi_n(1) &= \frac{\Gamma(M-n)}{\Gamma(M)} \sum_{i=0}^{n-1} (-1)^i \binom{n-1}{i} \\ &\quad \times ((1 - i\delta)^{M-1} - [1 - (i+1)\delta]^{M-1}) \\ &= \frac{\Gamma(M-n)}{\Gamma(M)} \sum_{i=0}^n \binom{n}{i} (-1)^i (1 - i\delta)^{M-1} \\ &= \frac{\Gamma(M-n)}{\Gamma(M)} \sum_{k=0}^{M-1} \binom{M-1}{k} (-1)^k \\ &\quad \times \left[ \sum_{i=0}^n \binom{n}{i} (-1)^i i^k \right] \delta^k. \end{aligned} \quad (81)$$

Since [27]

$$\sum_{k=0}^N \binom{N}{k} (-1)^k k^{(n-1)} = 0, \quad 1 \leq n \leq N, \quad (82)$$

$$\sum_{k=0}^N \binom{N}{k} (-1)^k k^N = (-1)^N N!, \quad N \geq 0, \quad (83)$$

we obtain  $\varphi_n(1) = \frac{\Gamma(M-n)}{\Gamma(M)} \sum_{k=n}^{M-1} \binom{M-1}{k} (-1)^k [\sum_{i=0}^n \binom{n}{i} (-1)^{i;k}] \delta^k$ . Substituting into (70) yields (27).

### APPENDIX C

#### PROOF OF Lemma 4

Similar to the proof of Lemma 3, we assume  $\mathbf{q}_i = \mathbf{e}_i$  without loss of generality. Then the numerator of (26) is given by

$$\begin{aligned} & \Pr \left( \sum_{i=1}^{n-1} |\mathbf{v}_k^H \mathbf{q}_i^H|^2 \leq 1-x, |\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta, \dots, |\mathbf{v}_k^H \mathbf{q}_{n-1}^H|^2 < \delta \right) \\ &= \Pr \left( \sum_{i=1}^{n-1} |v_i|^2 \leq 1-x, |v_1|^2 < \delta, \dots, |v_{n-1}|^2 < \delta \right). \end{aligned} \quad (84)$$

Recalling that  $t_i = |v_i|^2$ ,  $i = 1, 2, \dots, n-1$ , we can evaluate (84) using the joint p.d.f.  $f(t_1, \dots, t_{n-1})$  given in (69) in Appendix B. For  $n = 2$ , we have

$$\Pr(|\mathbf{v}_k^H \mathbf{q}_1^H|^2 \leq 1-x, |\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta) = \begin{cases} \int_0^\delta (M-1)(1-t_1)^{M-2} dt_1 & x \leq 1-\delta \\ \int_0^{1-x} (M-1)(1-t_1)^{M-2} dt_1 & 1-\delta < x \leq 1 \\ 0 & x > 1 \end{cases} \quad (85)$$

Solving the integrals in (85) and combining the result with (27) and (26) leads to the explicit solution given in (31). For  $n > 2$ , the problem is much more difficult. In this case, using (69), we obtain

$$\begin{aligned} & \Pr \left( \sum_{i=1}^{n-1} |\mathbf{v}_k^H \mathbf{q}_i^H|^2 \leq 1-x, |\mathbf{v}_k^H \mathbf{q}_1^H|^2 < \delta, \dots, |\mathbf{v}_k^H \mathbf{q}_{n-1}^H|^2 < \delta \right) \\ &= \begin{cases} 0 & x > 1 \\ \mu_n(\delta) & x \leq 1 - (n-1)\delta \\ \frac{\Gamma(M)}{\Gamma(M-n+1)} \int_{t_{n-1}} \dots \int_{t_1} \left(1 - \sum_{i=1}^{n-1} t_i\right)^{M-n} dt_1 \dots dt_{n-1} & 1 - (n-1)\delta < x \leq 1 \end{cases} \end{aligned} \quad (86)$$

with the integration region for the remaining multi-dimensional integral defined in the lemma statement.

Combining (86) with (27) and (26) leads to (30).

## APPENDIX D

## PROOF OF Lemma 5

We can upper bound the c.d.f. (30), for  $n \geq 2, 1 - (n - 1)\delta < x \leq 1$ , as follows

$$\begin{aligned}
F_{\beta(n)}(x) &\leq 1 - \frac{\Gamma(M)}{\Gamma(M - n + 1)\mu_n(\delta)} \\
&\times \int_0^{\frac{1-x}{n-1}} \cdots \int_0^{\frac{1-x}{n-1}} \left(1 - \sum_{i=1}^{n-1} t_i\right)^{M-n} dt_1 \cdots dt_{n-1} \\
&= 1 - \frac{\mu_n\left(\frac{1-x}{n-1}\right)}{\mu_n(\delta)}
\end{aligned} \tag{87}$$

where the second line follows from (70). For  $n = 2$ , we have

$$F_{\beta(2)}(x) \leq 1 - \frac{\mu_2(1-x)}{\mu_2(\delta)} = \frac{x^{M-1} - (1-\delta)^{M-1}}{(1-\delta)^{M-1}} \tag{88}$$

which is exactly the right-hand side of (31).

We can establish the corresponding lower bound via

$$\begin{aligned}
F_{\beta(n)}(x) &\geq 1 - \frac{\Gamma(M)}{\Gamma(M - n + 1)\mu_n(\delta)} \\
&\times \int_{\substack{\sum_{i=1}^{n-1} t_i \leq 1-x \\ t_1 \geq 0, \dots, t_{n-1} \geq 0}} \cdots \int \left(1 - \sum_{i=1}^{n-1} t_i\right)^{M-n} dt_1 \cdots dt_{n-1} \\
&= 1 - \frac{\Gamma(M)}{\Gamma(M - n + 1)\mu_n(\delta)} \\
&\times \int_0^{1-x} (1-y)^{M-n} \frac{y^{n-2}}{(n-2)!} dy \\
&= 1 - \frac{I_{1-x}(n-1, M-n+1)}{\mu_n(\delta)},
\end{aligned} \tag{89}$$

where we have used the identity [27]  $\int \cdots \int_{\substack{\sum_{i=1}^n t_i \leq h \\ t_1 \geq 0, \dots, t_n \geq 0}} dt_1 \cdots dt_n = \frac{h^n}{n!}$ . For  $n = 2$ , it is easily verified that (89) is equal to (88).

## APPENDIX E

## PROOF OF Lemma 6

Recalling that for uncorrelated Wishart matrices, the eigenvalues and their corresponding eigenvectors are independent, it follows that  $\lambda_{k,\max}$  is independent of  $\beta_k(n)$ ,  $\tilde{\beta}_k(n)$ , and  $\bar{\beta}_k(n)$ . Thus, the c.d.f.s of  $\gamma_k(n)$ ,  $\tilde{\gamma}_k(n)$ , and  $\bar{\gamma}_k(n)$ , can be derived as  $F_{\gamma(n)}(x) = \int_0^\infty F_{\beta(n)}(x/y) f_{\max}(y) dy$ ,  $F_{\tilde{\gamma}(n)}(x) = \int_0^\infty F_{\tilde{\beta}(n)}(x/y) f_{\max}(y) dy$ , and  $F_{\bar{\gamma}(n)}(x) = \int_0^\infty F_{\bar{\beta}(n)}(x/y) f_{\max}(y) dy$  respectively, where  $f_{\max}(\cdot)$  is the

p.d.f. of the maximum eigenvalue of  $\mathbf{H}_k \mathbf{H}_k^H$ . Together with *Lemma 5*, it follows trivially that  $F_{\tilde{\gamma}(n)}(x) \leq F_{\gamma(n)}(x) \leq F_{\tilde{\gamma}(n)}(x)$ , where the equalities hold for  $n = 2$ .

What remains is to derive closed-form expressions for  $F_{\tilde{\gamma}(n)}(x)$  and  $F_{\tilde{\beta}(n)}(x)$ . First consider  $F_{\tilde{\gamma}(n)}(x)$ . Recalling (32), and noting that for  $1 - (n-1)\delta < x \leq 1$ ,  $F_{\tilde{\beta}(n)}(x)$  can be re-expressed using (27) as

$$\begin{aligned} F_{\tilde{\beta}(n)}(x) &= 1 - \frac{1}{\mu_n(\delta)} \sum_{k=n-1}^{M-1} \binom{M-1}{k} (-1)^k \\ &\quad \times \left[ \sum_{i=0}^{n-1} \binom{n-1}{i} (-1)^i \left( \frac{i}{n-1} \right)^k (1-x)^k \right] \end{aligned} \quad (90)$$

it follows using *Lemma 1* that

$$\begin{aligned} F_{\tilde{\gamma}(n)}(x) &= F_{\max} \left( \frac{x}{t} \right) - \frac{1}{\mu_n(\delta)} \sum_{k=n-1}^{M-1} \binom{M-1}{k} (-1)^k \\ &\quad \times \left[ \sum_{i=0}^{n-1} \binom{n-1}{i} (-1)^i \left( \frac{i}{n-1} \right)^k \right] \sum_{r=1}^p \sum_{s=q-p}^{(N+M-2r)r} \\ &\quad a_{r,s} \int_x^{\frac{x}{t}} \left( 1 - \frac{x}{y} \right)^k y^s e^{-ry} dy. \end{aligned} \quad (91)$$

By applying the transformation  $z = \frac{y}{x}$  along with some elementary algebraic manipulations, the remaining integral is evaluated as

$$\begin{aligned} &\int_x^{\frac{x}{t}} \left( 1 - \frac{x}{y} \right)^k \frac{y^s}{e^{ry}} dy \\ &= \sum_{j=0}^k \binom{k}{j} (-1)^{k-j} r^{k-j-s-1} x^{k-j} \\ &\quad \times \left[ \Gamma(j-k+s+1, rx) - \Gamma \left( j-k+s+1, \frac{rx}{t} \right) \right]. \end{aligned}$$

Substituting this expression into (91), we readily obtain the result (34). A closed-form expression for  $F_{\tilde{\gamma}(n)}(x)$  can be obtained in a similar manner, and is omitted due to space limitations.

## APPENDIX F

### ASYMPTOTIC EXPANSION OF C.D.F.S OF $\tilde{\gamma}_k(n)$ AND $\bar{\gamma}_k(n)$ FOR LARGE $x$

First note that the tail behavior (large  $x$ ) of  $F_{\max}(x)$  is given by [15]

$$F_{\max}(x) = 1 - \frac{e^{-x} x^{M+N-2}}{\Gamma(M)\Gamma(N)} + O(e^{-x} x^{M+N-3}). \quad (92)$$

Then, the corresponding expansion for the term  $F_{\max}(\frac{x}{t})$  in both (34) and (35) follows immediately. In the following, we require a corresponding expansion for the remaining terms in (34) and (35). First consider (34). Since the remaining terms in this case involve the upper incomplete gamma function  $\Gamma(n, x)$ , we

require an asymptotic expansion for  $\Gamma(n, x)$  at  $x \rightarrow \infty$ . Using the definition and integrating by parts, for large  $x$  we have  $\Gamma(n, x) = e^{-x} x^{n-1} [1 + \frac{n-1}{x} + \frac{(n-1)(n-2)}{x^2} + \dots]$ . Since  $t < 1$ , the terms that decay most slowly in the summation in (34) can be expressed as

$$\begin{aligned} \mathcal{J}_1 &= \sum_{k=n-1}^{M-1} \mathcal{C}_k \sum_{s=q-p}^{N+M-2} \frac{a_{1,s} x^s}{e^x} \sum_{j=0}^k \frac{\binom{k}{j}}{(-1)^{k-j}} \\ &\quad \times \left[ 1 + \frac{j-k+s}{x} + \frac{(j-k+s)(j-k+s-1)}{x^2} + \dots \right], \end{aligned} \quad (93)$$

where

$$\mathcal{C}_k = \binom{M-1}{k} (-1)^k \left[ \sum_{i=0}^{n-1} \binom{n-1}{i} (-1)^i \left( \frac{i}{n-1} \right)^k \right]. \quad (94)$$

Using (82) we can obtain

$$\begin{aligned} \sum_{j=0}^k \binom{k}{j} (-1)^{k-j} j^{(m-1)} &= 0, \quad 1 \leq m \leq k, \\ \sum_{j=0}^k \binom{k}{j} (-1)^{k-j} j^k &= k!, \quad k \geq 1, \end{aligned} \quad (95)$$

from which it follows that in (93),  $\sum_{j=0}^k \binom{k}{j} (-1)^{k-j} \frac{\prod_{v=1}^m (j-k+s+1-v)}{x^m} = 0$  for  $1 \leq m < k-1$ , and also that  $\sum_{j=0}^k \binom{k}{j} (-1)^{k-j} \frac{\prod_{v=1}^k (j-k+s+1-v)}{x^k} = \frac{k!}{x^k}$ . We then have

$$\mathcal{J}_1 = \sum_{k=n-1}^{M-1} \mathcal{C}_k \sum_{s=q-p}^{N+M-2} \frac{a_{1,s} x^s k!}{e^x} \left( \frac{1}{x^k} + O\left(\frac{1}{x^{k+1}}\right) \right), \quad (96)$$

which upon substituting for  $\mathcal{C}_k$  and applying some manipulations using (95) gives

$$\begin{aligned} \mathcal{J}_1 &= \frac{(M-1)!(n-1)!}{(M-n)!(n-1)^{n-1}} a_{1,M+N-2} e^{-x} x^{M+N-n-1} \\ &\quad + O(e^{-x} x^{M+N-n-2}). \end{aligned} \quad (97)$$

From (92), we have  $f_{\max}(x) = \frac{e^{-x} x^{N+M-2}}{\Gamma(M)\Gamma(N)} + O(e^{-x} x^{N+M-3})$ . Therefore  $a_{1,N+M-2} = \frac{1}{\Gamma(M)\Gamma(N)}$ . Together with (97) and (92), we have (36). By using a similar method, the terms that decay most slowly in the summation in (35) can be obtained. That result, used with (92), yields (37).

## APPENDIX G

### PROOF OF Lemma 8

Recall that  $F_{\tilde{\gamma}_{\pi(n)}}(x) \leq F_{\gamma_{\pi(n)}}(x) \leq F_{\tilde{\gamma}_{\pi(n)}}(x)$ . For  $\gamma_{\pi(n)}$ ,  $n \in \{2, \dots, M\}$ , and large  $K$ , with (40),  $\Pr\{u_n - \log \log \sqrt{K} \leq \gamma_{\pi(n)}\} \geq \Pr\{u_n - \log \log \sqrt{K} \leq \tilde{\gamma}_{\pi(n)}\} \geq 1 - O\left(\frac{1}{\log K}\right)$ . Similarly, with (41)

we have  $\Pr\{\gamma_{\pi(n)} \leq \chi_n + \log \log \sqrt{K}\} \geq \Pr\{\bar{\gamma}_{\pi(n)} \leq \chi_n + \log \log \sqrt{K}\} \geq 1 - O\left(\frac{1}{\log K}\right)$ . Thus,

$$\begin{aligned} & \Pr\{u_n - \log \log \sqrt{K} \leq \gamma_{\pi(n)} \leq \chi_n + \log \log \sqrt{K}\} \\ & \geq 1 - O\left(\frac{1}{\log K}\right). \end{aligned} \quad (98)$$

For  $n = 1$ , the asymptotic distribution of  $\gamma_{\pi(n)}$  has been characterized in [14]. Using that result, along with (98), the lemma follows upon noting that  $\zeta_{\pi(n)} = \rho\gamma_{\pi(n)}$ .

## APPENDIX H

### PROOF OF *Theorem 1*

Using (45) we can obtain  $\Pr\left\{\frac{\log_2(1+\varpi_n-\rho\log\log\sqrt{K})}{\log_2[\rho\log K]} \leq \frac{\log_2(1+\zeta_{\pi(n)})}{\log_2[\rho\log K]} \leq \frac{\log_2(1+v_n+\rho\log\log\sqrt{K})}{\log_2[\rho\log K]}\right\} \geq 1 - O\left(\frac{1}{\log K}\right)$ . Substituting (46) and (47) and letting  $K \rightarrow \infty$ , the left-hand side and right-hand side inequality within  $\Pr\{\cdot\}$  converge to the same value. Thus,  $\lim_{K \rightarrow \infty} \frac{\log_2(1+\zeta_{\pi(n)})}{\log_2[\rho\log K]} = 1$  with probability 1, and (48) holds. To establish (49), we employ the following upper bound on  $\mathcal{E}\{R_{\text{BC}}\}$  derived in [16]:

$$\mathcal{E}\{R_{\text{BC}}\} \leq M \log_2(1 + \rho(\log K + O(\log \log K))). \quad (99)$$

From *Lemma 8*, we have  $\Pr\left\{\log_2(1 + \zeta_{\pi(n)}) \geq \log_2(1 + \varpi_n - \rho \log \log \sqrt{K})\right\} \geq 1 - O\left(\frac{1}{\log K}\right)$ . Thus,

$$\begin{aligned} & \mathcal{E}\{R_{\text{BC}}\} - \mathcal{E}\{R_{\text{ZFDP-C-SUS}}\} \\ & \leq M \log(1 + \rho(\log K + O(\log \log K))) \\ & \quad - \left(1 - O\left(\frac{1}{\log K}\right)\right) \\ & \quad \times \sum_{n=1}^M \log(1 + \varpi_n - \rho \log \log \sqrt{K}) \\ & \sim \sum_{n=1}^M \log\left(1 + \frac{O(\log \log K)}{1 + \varpi_n - \rho \log \log \sqrt{K}}\right) \\ & \quad + O\left(\frac{1}{\log K}\right) M O(\log \log K) \\ & \sim O\left(\frac{\log \log K}{\log K}\right) \end{aligned} \quad (100)$$

where we have used  $\log(1+x) \approx x$  for  $x \ll 1$ , and  $x \sim y$  means  $\lim_{K \rightarrow \infty} x/y = 1$ .



## APPENDIX I

PROOF OF *Theorem 2*

From [13], for small enough  $\delta$ ,  $\varrho_{\pi(n)} > \frac{\gamma_{\pi(n)}}{1+e(\delta)}$ , where  $e(\delta) = \frac{(M-1)^4\delta}{1-(M-1)\delta}$ . Using this result, together with (99) and (45), and following a similar method as in Appendix H, we have

$$\begin{aligned}
\mathcal{E}\{R_{\text{BC}}\} & - \mathcal{E}\{R_{\text{ZFBF-SUS}}\} \\
& \leq M \log \left( 1 + \rho(\log K + O(\log \log K)) \right) - \mathcal{E} \left\{ \sum_{n=1}^M \log \left( 1 + \frac{\rho \gamma_{\pi(n)}}{1 + e(\delta)} \right) \right\} \\
& \leq M \log \left( 1 + \rho(\log K + O(\log \log K)) \right) - \sum_{n=1}^M \left( 1 - O \left( \frac{1}{\log K} \right) \right) \log \left( 1 + \frac{\varpi_n - \rho \log \log \sqrt{K}}{1 + e(\delta)} \right) \\
& \sim \sum_{n=1}^M \log \left( 1 + \frac{\rho(e(\delta) \log K + O(\log \log K))}{1 + (\varpi_n - \rho \log \log \sqrt{K}) \sum_{i=0}^{\infty} (-e(\delta))^i} \right) + O \left( \frac{\log \log K}{\log K} \right) \\
& \sim M e(\delta) + O \left( \frac{\log \log K}{\log K} \right), \tag{101}
\end{aligned}$$

where we have used the fact that for small enough  $\delta$ ,  $|e(\delta)| < 1$ , thus  $\frac{1}{1+e(\delta)} = \sum_{i=0}^{\infty} (-e(\delta))^i$ . So we can see that as long as  $e(\delta) \sim o(1)$ , or equivalently  $\delta \sim o(1)$ , whilst satisfying the conditions in (29), the difference will become zero as  $K \rightarrow \infty$ . However, obviously ZFBF-SUS with a smaller candidate set at each iteration (i.e., reduced  $|\mathcal{U}_n|$ ) can not achieve more sum rate than ZFBF-SUS with a larger candidate set at each iteration. Thus, with larger  $\delta$ , there will be more candidate users for each iteration and the average sum rate will increase, or at least maintain. So the condition  $\delta \sim o(1)$  can be ignored, thereby establishing (50). From (101), the difference in sum rate is at most  $O\left(\frac{\log \log K}{\log K}\right)$ .

## REFERENCES

- [1] W. Yu and J. M. Cioffi, "Sum capacity of a Gaussian vector broadcast channels," *IEEE Trans. Inform. Theory*, vol. 50, no. 9, pp. 1875–1892, Sep. 2002.
- [2] P. Viswanath and D. N. C. Tse, "Sum capacity of the vector Gaussian broadcast channel and uplink-downlink duality," *IEEE Trans. Inform. Theory*, vol. 49, no. 8, pp. 1912–1921, Aug. 2003.
- [3] H. Weingarten, Y. Steinberg, and S. Shamai (Shitz), "The capacity region of the Gaussian multiple-input multiple-output broadcast channel," *IEEE Trans. Inform. Theory*, vol. 52, no. 9, pp. 3936–3964, Sep. 2006.
- [4] S. Vishwanath, N. Jindal, and A. Goldsmith, "Duality, achievable rates, and sum-rate capacity of Gaussian MIMO broadcast channels," *IEEE Trans. Inform. Theory*, vol. 49, no. 10, pp. 2658–2668, Oct. 2003.
- [5] G. Caire and S. Shamai (Shitz), "On the achievable throughput of a multi-antenna Gaussian broadcast channel," *IEEE Trans. Inform. Theory*, vol. 49, no. 7, pp. 1691–1706, Jul. 2003.
- [6] A. D. Dabbagh and D. J. Love, "Precoding for multiple antenna Gaussian broadcast channels with successive zero-forcing," *IEEE Trans. Signal Process.*, vol. 55, no. 7, pp. 3837–3850, Jul. 2007.
- [7] C. B. Peel, B. M. Hochwald, and A. L. Swindlehurst, "A vector-perturbation technique for near-capacity multiantenna multi-user communication - Part I: Channel inversion and regularization," *IEEE Trans. Commun.*, vol. 53, no. 1, pp. 195–202, Jan. 2005.

- [8] B. M. Hochwald, C. B. Peel, and A. L. Swindlehurst, "A vector-perturbation technique for near-capacity multi-antenna multiuser communication - Part II: Perturbation," *IEEE Trans. Commun.*, vol. 53, no. 3, pp. 537–544, Mar. 2005.
- [9] Q. H. Spencer, A. L. Swindlehurst, and M. Haardt, "Zero-forcing methods for downlink spatial multiplexing in multiuser MIMO channels," *IEEE Trans. Signal Process.*, vol. 52, no. 2, pp. 461–471, Feb. 2004.
- [10] Z. Tu and R. S. Blum, "Multiuser diversity for a dirty paper approach," *IEEE Commun. Lett.*, vol. 7, no. 8, pp. 370–372, Aug. 2003.
- [11] G. Dimic and N. Sidiropoulos, "On the downlink beamforming with greedy user selection: Performance analysis and a simple new algorithm," *IEEE Trans. Signal Process.*, vol. 53, no. 10, pp. 3857–3868, Jul. 2005.
- [12] Z. Shen, R. Chen, J. G. Andrews, R. W. Heath Jr., and B. L. Evans, "Low complexity user selection algorithms for multiuser MIMO systems with block diagonalization," *IEEE Trans. Signal Process.*, vol. 54, no. 9, pp. 3658–3663, Sep. 2006.
- [13] T. Yoo and A. J. Goldsmith, "On the optimality of multi-antenna broadcast scheduling using zero-forcing beamforming," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 3, pp. 528–541, Mar. 2006.
- [14] M. A. Maddah-Ali, M. Ansari, and A. K. Khandani, "Broadcast in MIMO systems based on a generalized QR decomposition: signaling and performance analysis," *IEEE Trans. Inform. Theory*, vol. 54, no. 3, pp. 1124–1138, Mar. 2008.
- [15] A. Bayesteh and A. K. Khandani, "On the user selection for MIMO broadcast channels," *IEEE Trans. Inform. Theory*, vol. 54, no. 3, pp. 1086–1107, Mar. 2008.
- [16] M. Sharif and B. Hassibi, "On the capacity of MIMO broadcast channels with partial side information," *IEEE Trans. Inform. Theory*, vol. 2, no. 21, pp. 506–522, Feb. 2005.
- [17] —, "A comparison of time-sharing, DPC, and beamforming for MIMO broadcast channels with many users," *IEEE Trans. Commun.*, vol. 55, no. 1, pp. 11–15, Jan. 2007.
- [18] J. Wang, D. J. Love, and M. D. Zoltowski, "User selection with zero-forcing beamforming achieves the asymptotically optimal sum rate," *IEEE Trans. Signal Process.*, vol. 56, no. 8, pp. 3713–3726, Aug. 2008.
- [19] H. Viswanathan, S. Venkatesan, and H. Huang, "Downlink capacity evaluation of cellular networks with known-interference cancellation," *IEEE J. Sel. Areas Commun.*, vol. 21, no. 6, pp. 802–811, Jun. 2003.
- [20] P. A. Dighe, R. K. Mallik, and S. S. Jamuar, "Analysis of transmit-receive diversity in Rayleigh fading," *IEEE Trans. Commun.*, vol. 51, no. 4, pp. 694–703, Apr. 2003.
- [21] A. Maaref and S. Aïssa, "Closed-form expressions for the outage and ergodic Shannon capacity of MIMO MRC systems," *IEEE Trans. Commun.*, vol. 53, no. 7, pp. 1092–1095, Jul. 2005.
- [22] H. David and H. Nagaraja, *Order Statistics*, 3rd ed. New York: John Wiley and Sons, 2003.
- [23] X. Zhang and S.-Y. Kung, "Capacity analysis for parallel and sequential MIMO equalizers," *IEEE Trans. Signal Process.*, vol. 11, no. 51, pp. 2989–3002, Nov. 2003.
- [24] N. Jindal, W. Rhee, S. Vishwanath, S. Jafar, and A. Goldsmith, "Sum power iterative water-filling for multi-antenna Gaussian broadcast channels," *IEEE Trans. Inform. Theory*, vol. 51, no. 4, pp. 1570–1580, Apr. 2005.
- [25] J. Wang, D. J. Love, and M. D. Zoltowski, A Result on Order Statistics. [Online]. Available: <http://docs.lib.purdue.edu/ecetr/347>, Tech. Rep., Purdue Univ., West Lafayette, IN, 2007.
- [26] M. G. Kendall, *A course in the geometry of n dimensions*, 1st ed. London, U.K.: Charles Griffin Co., Ltd., 1961.
- [27] I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products*, 6th ed. New York: Academic, 2000.