# Statistical Beamforming on the Grassmann Manifold for the Two-User Broadcast Channel

Vasanthan Raghavan\*, Stephen V. Hanly, Venugopal V. Veeravalli

#### Abstract

A Rayleigh fading spatially correlated broadcast setting with M = 2 antennas at the transmitter and two-users (each with a single antenna) is considered. It is assumed that the users have perfect channel information about their links whereas the transmitter has only statistical information of each user's link (covariance matrix of the vector channel). A low-complexity linear beamforming strategy that allocates equal power and one spatial eigen-mode to each user is employed at the transmitter. Beamforming vectors on the Grassmann manifold that depend only on statistical information are to be designed at the transmitter to maximize the ergodic sum-rate delivered to the two users. Towards this goal, the beamforming vectors are first fixed and a closed-form expression is obtained for the ergodic sum-rate in terms of the covariance matrices of the links. This expression is non-convex in the beamforming vectors ensuring that the classical Lagrange multiplier technique is not applicable. Despite this difficulty, the optimal solution to this problem is shown to be the solution to the maximization of an appropriatelydefined average signal-to-interference and noise ratio (SINR) metric for each user. This solution is the dominant generalized eigenvector of a pair of positive-definite matrices where the first matrix is the covariance matrix of the forward link and the second is an appropriately-designed "effective" interference covariance matrix. In this sense, our work is a generalization of optimal signalling along the dominant eigen-mode of the transmit covariance matrix in the single-user case. Finally, the ergodic sum-rate

V. Raghavan is with the Department of Mathematics, University of Southern California, Los Angeles, CA 90089, USA. He was with the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA and the Department of Electrical and Electronic Engineering, The University of Melbourne, Parkville, VIC 3052, Australia when parts of this work was done. S. V. Hanly is with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore 117576. V. V. Veeravalli is with the Coordinated Science Laboratory and the Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA. Email: vasanthan\_raghavan@ieee.org, elehsv@nus.edu.sg, vvv@illinois.edu. \*Corresponding author.

This work has been supported in part by the ARC through grant DP-0984862, the NUS through grant WBS #R-263-000-572-133 and the NSF through grant CNS-0831670 at the University of Illinois. This paper was presented in part at the IEEE International Symposium on Information Theory, Austin, TX, 2010. for the general broadcast setting with M antennas at the transmitter and M-users (each with a single antenna) is obtained in terms of the covariance matrices of the links and the beamforming vectors.

#### **Index Terms**

Adaptive signalling, broadcast channel, information rates, MISO systems, multi-user MIMO, precoding, spatial correlation.

#### I. INTRODUCTION

The last fifteen years of research in wireless communications has seen the emergence of multi-antenna signalling as a viable option to realize high data-rates at practically acceptable reliability levels. While initial work on multi-antenna design was primarily motivated by the single-user paradigm [1]–[5], more recent attention has been on the theory and practice of multi-user multi-antenna communications [6]–[9]. The focus of this paper is on a broadcast setting that typically models a cellular downlink. We study the multiple-input single-output (MISO) broadcast problem where a central transmitter with M antennas communicates with M users in the cell, each having a single antenna. Under the assumption of perfect channel state information (CSI) at both the transmitter and the user ends, significant progress has been made over the last few years on understanding optimal signalling that achieves the sum-capacity [10]–[15] as well as the capacity region [16] of the multi-antenna broadcast channel. The capacity-achieving *dirty-paper coding* scheme [17] pre-nulls interference from simultaneous transmissions by other users to a specific user and hence results in a multiplexing gain of M.

Nevertheless, the high implementation complexity associated with dirty-paper coding [18] makes it less attractive in standardization efforts for practical systems. The consequent search for low-complexity signalling alternatives that are within a fixed power-offset<sup>1</sup> of the dirty-paper coding scheme has resulted in an array of candidate linear (as well as non-linear) precoding techniques [19]–[27]. In particular, a linear beamforming scheme that is developed as a generalization of the single-user beamforming scheme has attracted significant attention in the literature. Specifically, a scheme where the transmitter allocates one eigen-mode to each user and shares the power budget equally among all the users is the focus of this work.

If perfect CSI is available at both the ends, instantaneous nulls can be created in the interference sub-space of each user (or interference can be zeroforced) and thus this scheme remains orderoptimal with respect to the dirty-paper coding scheme [23]. However, the practical utility of the linear beamforming scheme is dependent on how gracefully its performance degrades with the

<sup>&</sup>lt;sup>1</sup>Two schemes are within a fixed power-offset if the difference in power level necessary to achieve a fixed rate with the two schemes stays bounded independent of the rate.

quality of CSI at the transmitter. This is because while reasonably accurate CSI can be obtained at the user end via pilot-based training schemes, CSI at the transmitter requires either channel reciprocity or reverse link feedback, both of which put an overwhelming burden on the operating cost [9]. In the extreme (and pessimistic) setting of no CSI at the transmitter, the multiplexing gain reduces to 1 (that is, it is lost completely relative to the perfect CSI case).

In practice, the channel evolves fairly slowly on a statistical scale and it is possible to learn the spatial statistics<sup>2</sup> of the individual links at the transmitter with minimal cost. With only statistical information at the transmitter, the interference cannot be nulled out completely and a low-complexity decoder architecture that treats interference as noise is often preferred. Initial works assume an identity covariance matrix for all the users corresponding to an *independent and identically distributed* (i.i.d.) fading process in the spatial domain [6]–[9]. However, this model cannot be justified in practical systems that are often deployed in environments where the scattering is localized in certain spatial directions or where antennas are not spaced wide apart due to infrastructural constraints [28].

While signalling design for the single-user setting under a very general spatial correlation model is now well-understood [1]–[5], the broadcast case where the channel statistics vary across users and different users experience different covariance matrices has not received much attention. In particular, [29] studies the problem where all the users share a common non-i.i.d. transmit covariance matrix and captures the impact of this common covariance matrix on the achievable rates. In [30], the authors show that second-order spatial statistics can be exploited to schedule users that enjoy better channel quality and hence improve the overall performance of an opportunistic beamforming scheme. In the same spirit, it is shown in [31] and [32] that second-order moments of the channel in combination with instantaneous norm (or weighted-norm) feedback is sufficient to extract almost all of the multi-user diversity gain in a broadcast setting. Spatial correlation is exploited to reduce the feedback overhead of a limited feedback codebook design in [33]–[36].

Summary of Main Contributions: With this background, the main focus of this paper is to fill some of the gaps in understanding the information-theoretic limits of broadcast channels with low-complexity signalling schemes (such as linear beamforming) under practical assumptions on CSI and decoder architecture. We study the simplest non-trivial version of this problem corresponding to the two-user (M = 2) case. We design optimal beamforming vectors on the Grassmann manifold<sup>3</sup>  $\mathcal{G}(2,1)$  to maximize the ergodic sum-rate achievable with the linear

 $<sup>^{2}</sup>$ With a Rayleigh (or a Ricean) fading model for the MISO channel, the complete statistical information of the link is captured by the covariance matrix (or the mean vector and the covariance matrix) of the vector channel.

<sup>&</sup>lt;sup>3</sup>Informally,  $\mathcal{G}(M, 1)$  denotes the space of all *M*-dimensional unit-norm beamforming vectors modulo the phase of the first element of the vector. A more formal definition is provided in Sec. II (Def. 1).

beamforming scheme.

The first step to this goal is the computation of the ergodic sum-rate in closed-form. For this, we develop insight into the structure of the density function of the weighted-norm of beamforming vectors isotropically distributed on  $\mathcal{G}(2,1)$ . Exploiting this knowledge, we derive an explicit expression for the ergodic sum-rate in terms of the covariance matrices of the users and the beamforming vectors. This expression can be rewritten in terms of a certain generalized "distance" measure between the beamforming vectors. As a result of this complicated non-linear dependence, the sum-rate is non-convex in the beamforming vectors thus precluding the use of the classical Lagrangian approach to convex optimization. Instead, a first-principles based technique is developed where the beamforming vectors are decomposed along an appropriately chosen (in general, non-orthogonal) basis. Exploiting this decomposition structure, we obtain an upper bound for the ergodic sum-rate, which we show is tight for a specific choice of beamforming vectors (see Theorems 2 and 3). This optimal choice is the dominant generalized eigenvector<sup>4</sup> of a pair of covariance matrices, with one of them being the covariance matrix of the forward link and the other an appropriately-designed "effective" interference covariance matrix. The generalized eigenvector structure is the solution to maximizing an appropriately-defined average signal-tointerference and noise ratio (SINR) metric for each user and thus generalizes our intuition from the single-user case [1]–[5]. Table I in the Conclusions section (Sec. VI) summarizes the structure of the optimal beamforming vectors under different signal-to-noise ratio (SNR) assumptions.

While a generalized eigenvector solution has been obtained in the perfect CSI case for the multiple-input multiple-output (MIMO) broadcast problem [26], [37] and the MIMO interference channel problem in the low-interference regime [38], to the best of our knowledge, its appearance in the statistical setting is a first. A closely-related work of ours [39] reports the optimality of the generalized eigenvector solution for the statistical beamformer design in the MISO interference channel setting with two antennas. We also extend our intuition to the weighted ergodic sumrate maximization problem [40] and conjecture on the structure of the optimal beamforming vectors. Numerical results justify our conjecture and the intuition behind it. Finally, closed-form expression for the ergodic sum-rate in terms of the covariance matrices of the links and the beamforming vectors are obtained in the general M-user case.

**Organization:** This paper is organized as follows. With Section II explaining the background of the problem, ergodic rate expressions in terms of the covariance matrices of the links and beamforming vectors are obtained in Section III for the M = 2 case. The non-convex optimization problem of ergodic sum-rate maximization is the main focus of Section IV with the low- and

<sup>&</sup>lt;sup>4</sup>A generalized eigenvector generalizes the notion of an eigenvector to a pair of matrices. A more technical definition is provided in Sec. IV (Def. 2). The dominant eigenvector is the eigenvector corresponding to the dominant eigenvalue. Under the assumption that the eigenvector is unit-norm, it is unique on  $\mathcal{G}(M, 1)$ .

the high-SNR extremes providing insight for the development in the intermediate-SNR regime. The focus shifts to weighted ergodic sum-rate maximization in Section V. In addition, sum-rate expressions are generalized to the general M-user case and concluding remarks are provided in Section VI. Most of the proofs/details are relegated to the Appendices.

*Notation:* We use upper- and lower-case bold symbols for matrices and vectors, respectively. The notations  $\Lambda$  and U are usually reserved for eigenvalue and eigenvector matrices whereas I is reserved for the identity matrix (of appropriate dimensionality). The *i*-th diagonal element of  $\Lambda$  is denoted by  $\Lambda_i$  while the *i*-th element of a vector  $\mathbf{x}$  is denoted by  $\mathbf{x}(i)$ . At times, we also use  $\lambda_1, \lambda_2, \cdots$  to denote the eigenvalues of a Hermitian matrix, and these eigenvalues are often arranged in decreasing order as  $\lambda_1 \geq \lambda_2 \geq \cdots$ . The Hermitian transpose and inverse operations of a matrix are denoted by  $(\cdot)^H$  and  $(\cdot)^{-1}$  while the trace operator is denoted by  $\mathrm{Tr}(\cdot)$ . The two-norm of a vector is denoted by the symbol  $\|\cdot\|$ . The operator  $E[\cdot]$  stands for expectation while the density function of a random variable is denoted by the symbol  $p(\cdot)$ . The symbols  $\mathbb{C}$  and  $\mathbb{R}^+$  stand for complex and positive real fields, respectively.  $X \sim \mathcal{CN}(\mu, \sigma^2)$  indicates that X is a complex Gaussian random variable with mean  $\mu$  and variance  $\sigma^2$ .

#### **II. SYSTEM SETUP**

We consider a broadcast setting that models a MISO cellular downlink with M antennas at the transmitter and M users, each with a single antenna. We denote the  $M \times 1$  vector channel between the transmitter and user i as  $\mathbf{h}_i$ ,  $i = 1, \dots, M$ . While different multi-user communication strategies can be considered [19]–[27], as motivated in Sec. I, the focus here is on a linear beamforming scheme where the information-bearing signal  $s_i$  meant for user iis beamformed from the transmitter with the  $M \times 1$  unit-norm vector  $\mathbf{w}_i$ . We assume that  $s_i$ is unit energy and the transmitter divides its power budget<sup>5</sup> of  $\rho$  equally across all the users. Equal power allocation is popular in current-generation cellular standards where low-complexity schemes are preferred. The received symbol  $y_i$  at user i is written as

$$y_i = \sqrt{\frac{\rho}{M}} \cdot \mathbf{h}_i^H \left( \sum_{i=1}^M \mathbf{w}_i s_i \right) + n_i, \ i = 1, \cdots, M$$
(1)

where  $\rho$  is the transmit power and  $n_i$  denotes the  $\mathcal{CN}(0,1)$  complex Gaussian noise added at the receiver.

Initial works on the broadcast problem assume that  $h_i$  is ergodic and it evolves over time and frequency in an i.i.d. fashion, and is spatially i.i.d. While the above assumption can be justified in the time and frequency axes with a frame-based and multi-carrier signalling approach (common

<sup>5</sup>The practically motivated power-control problem where different powers could be allocated to the different users is a related problem, but it is not studied here.

in current-generation systems), it cannot be justified along the spatial axis. This is because the channel variation in the spatial (antenna) domain cannot be i.i.d. unless the antennas at the transmitter end are spaced wide apart and the scattering environment connecting the transmitter with the users is sufficiently rich [28]. With this motivation, the main emphasis of this work is on understanding the impact of the users' spatial statistics on the performance of a linear beamforming scheme.

We assume a Rayleigh fading<sup>6</sup> (zero mean complex Gaussian) model for the channel, which implies that the complete spatial statistics are described by the second-order moments of  $\{\mathbf{h}_i\}$ . For the MISO model, the channel  $\mathbf{h}_i$  of user *i* can be written as

$$\mathbf{h}_i = \boldsymbol{\Sigma}_i^{1/2} \mathbf{h}_{\mathsf{iid}, i} \tag{2}$$

where  $\mathbf{h}_{iid,i}$  is an  $M \times 1$  vector with i.i.d.  $\mathcal{CN}(0,1)$  entries and  $\Sigma_i \triangleq E\left[\mathbf{h}_i \mathbf{h}_i^H\right]$  is the transmit covariance matrix corresponding to user *i*. Note that (2) is the most general statistical model for  $\mathbf{h}_i$  under the MISO assumption. With  $\Sigma_i = \mathbf{I}$  for all users, (2) reduces to the i.i.d. downlink model well-studied in the literature [6]–[9].

Under the assumption of Gaussian inputs  $\{s_i\}$ , the instantaneous information-theoretic<sup>7</sup> rate,  $R_i$ , achievable by user *i* with the linear beamforming scheme using a mismatched<sup>8</sup> decoder [41] is given by

$$R_i = \log\left(1 + \frac{\frac{\rho}{M} \cdot |\mathbf{h}_i^H \mathbf{w}_i|^2}{1 + \frac{\rho}{M} \cdot \sum_{j \neq i} |\mathbf{h}_i^H \mathbf{w}_j|^2}\right)$$
(3)

$$= \underbrace{\log\left(1 + \frac{\rho}{M} \cdot \sum_{j=1}^{M} |\mathbf{h}_{i}^{H} \mathbf{w}_{j}|^{2}\right)}_{I_{i,1}} - \underbrace{\log\left(1 + \frac{\rho}{M} \cdot \sum_{j \neq i} |\mathbf{h}_{i}^{H} \mathbf{w}_{j}|^{2}\right)}_{I_{i,2}}.$$
(4)

With the spatial correlation model assumed in (2), we can write  $I_{i,1}$  as

$$I_{i,1} = \log\left(1 + \frac{\rho}{M} \cdot \mathbf{h}_{\mathsf{iid},i}^{H} \boldsymbol{\Sigma}_{i}^{1/2} \left(\sum_{j=1}^{M} \mathbf{w}_{j} \mathbf{w}_{j}^{H}\right) \boldsymbol{\Sigma}_{i}^{1/2} \mathbf{h}_{\mathsf{iid},i}\right)$$
(5)

$$= \log\left(1 + \frac{\rho}{M} \cdot \mathbf{h}_{\mathsf{iid},\,i}^{H} \mathbf{V}_{i} \mathbf{\Lambda}_{i} \mathbf{V}_{i}^{H} \mathbf{h}_{\mathsf{iid},\,i}\right),\tag{6}$$

where we have used the following eigen-decomposition in (6):

$$\mathbf{V}_{i} \mathbf{\Lambda}_{i} \mathbf{V}_{i}^{H} = \mathbf{\Sigma}_{i}^{1/2} \left( \sum_{j=1}^{M} \mathbf{w}_{j} \mathbf{w}_{j}^{H} \right) \mathbf{\Sigma}_{i}^{1/2}$$
(7)

$$\Lambda_{i} = \operatorname{diag}([\Lambda_{i,1}, \cdots, \Lambda_{i,M}]), \quad \Lambda_{i,1} \geq \cdots \geq \Lambda_{i,M} \geq 0.$$
(8)

<sup>6</sup>While more general fading models such as Ricean or Nakagami-*m* models can be considered, this paper focuses on the Rayleigh model alone.

<sup>7</sup>All logarithms are to base e and all rate quantities are assumed to be in nats/s/Hz in this work.

<sup>8</sup>Here, the decoding rule is different from the optimal decoding rule due to the presence of multi-user interference.

Similarly, we can write  $I_{i,2}$  as

$$I_{i,2} = \log\left(1 + \frac{\rho}{M} \cdot \mathbf{h}_{\mathsf{iid},i}^{H} \widetilde{\mathbf{V}}_{i} \widetilde{\mathbf{\Lambda}}_{i} \widetilde{\mathbf{V}}_{i}^{H} \mathbf{h}_{\mathsf{iid},i}\right)$$
(9)

$$\widetilde{\mathbf{V}}_{i}\widetilde{\mathbf{\Lambda}}_{i}\widetilde{\mathbf{V}}_{i}^{H} = \Sigma_{i}^{1/2} \left(\sum_{j\neq i} \mathbf{w}_{j}\mathbf{w}_{j}^{H}\right) \Sigma_{i}^{1/2}$$
(10)

$$\widetilde{\Lambda}_{i} = \operatorname{diag}([\widetilde{\Lambda}_{i,1}, \cdots, \widetilde{\Lambda}_{i,M}]), \quad \widetilde{\Lambda}_{i,1} \geq \cdots \geq \widetilde{\Lambda}_{i,M} \geq 0.$$
(11)

The goal of this work is to maximize the throughput conveyed from the transmitter to the users by the choice of beamforming vectors. Specifically, the metric of interest is the ergodic sum-rate,  $\mathcal{R}_{sum}$ , achievable with the linear beamforming scheme:

$$\mathcal{R}_{\mathsf{sum}} \triangleq \sum_{i=1}^{M} E\left[R_i\right]. \tag{12}$$

For this, note that the achievable rate in (3) is invariant to transformations of the form  $\mathbf{w}_i \mapsto e^{j\theta} \mathbf{w}_i$ for any  $\theta$ . Coupled with the unit-norm assumption for  $\mathbf{w}_i$ , the space over which optimization is performed is precisely defined as follows.

Definition 1 (Stiefel and Grassmann Manifolds [42]): The uni-dimensional complex Stiefel manifold St(M, 1) refers to the unit-radius complex sphere in M-dimensions and is defined as

$$\mathsf{St}(M,1) = \left\{ \mathbf{x} \in \mathbb{C}^M : \|\mathbf{x}\| = 1 \right\}.$$
(13)

The uni-dimensional complex Grassmann manifold  $\mathcal{G}(M, 1)$  consists of the set of one-dimensional subspaces of St(M, 1). Here, a transformation of the form  $\mathbf{x} \mapsto e^{j\theta}\mathbf{x}$  (for any  $\theta$ ) is treated as invariant by considering all vectors of the form  $e^{j\theta}\mathbf{x}$  (for some  $\theta$ ) to belong to the one-dimensional sub-space spanned by  $\mathbf{x}$ .

The optimization objective is then to understand the structure of the beamforming vectors,  $\{\mathbf{w}_{i, \text{ opt}}\}$ , that maximize  $\mathcal{R}_{sum}$ :

$$\mathbf{w}_{i, \mathsf{opt}} = \underset{\mathbf{w}_i \in \mathcal{G}(M, 1)}{\arg \max} \mathcal{R}_{\mathsf{sum}}, \ i = 1, \cdots, M.$$
(14)

In (14), the candidate beamforming vectors,  $\{\mathbf{w}_i\}$ , depend only on the long-term statistics of the channel, which (as noted before) in the MISO setting is the set of all transmit covariance matrices,  $\{\Sigma_i\}$ .

Towards the goal of computing  $\mathcal{R}_{sum}$ , we decompose  $\mathbf{h}_{iid, i}$  into its magnitude and directional components as  $\mathbf{h}_{iid, i} = \|\mathbf{h}_{iid, i}\| \cdot \widehat{\mathbf{h}}_{iid, i}$ . It is well-known [1] that  $\|\mathbf{h}_{iid, i}\|^2$  can be written as

~ 1 (

$$\|\mathbf{h}_{\mathsf{iid},i}\|^2 = \frac{1}{2} \sum_{j=1}^{2M} \chi_j^2 \tag{15}$$

where  $\chi_j^2$  is a standard (real) chi-squared random variable and  $\widehat{\mathbf{h}}_{\text{iid},i}$  is a unit-norm vector that is isotropically distributed [42], [43] on  $\mathcal{G}(M, 1)$ . Thus, we can rewrite  $I_{i,1}$  and  $I_{i,2}$  as

$$I_{i,1} = \log\left(1 + \frac{\rho}{M} \cdot \|\mathbf{h}_{\mathsf{iid},i}\|^2 \cdot \widehat{\mathbf{h}}_{\mathsf{iid},i}^H \mathbf{V}_i \mathbf{\Lambda}_i \mathbf{V}_i^H \widehat{\mathbf{h}}_{\mathsf{iid},i}\right)$$
(16)

$$I_{i,2} = \log\left(1 + \frac{\rho}{M} \cdot \|\mathbf{h}_{\mathsf{iid},i}\|^2 \cdot \widehat{\mathbf{h}}_{\mathsf{iid},i}^H \widetilde{\mathbf{V}}_i \widetilde{\mathbf{\Lambda}}_i \widetilde{\mathbf{V}}_i^H \widehat{\mathbf{h}}_{\mathsf{iid},i}\right).$$
(17)

Since the magnitude and directional information of an i.i.d. random vector are independent [43],  $E[I_{i,1}]$  and  $E[I_{i,2}]$  can be further written as

$$E[I_{i,1}] = E_{\|\mathbf{h}_{\mathsf{iid},i}\|} \left[ E_{\widehat{\mathbf{h}}_{\mathsf{iid},i}} \left[ \log \left( 1 + \frac{\rho}{M} \cdot \|\mathbf{h}_{\mathsf{iid},i}\|^2 \cdot \widehat{\mathbf{h}}_{\mathsf{iid},i}^H \mathbf{\Lambda}_i \widehat{\mathbf{h}}_{\mathsf{iid},i} \right) \right] \right]$$
(18)

$$E\left[I_{i,2}\right] = E_{\|\mathbf{h}_{\mathsf{iid},\,i}\|} \left[ E_{\widehat{\mathbf{h}}_{\mathsf{iid},\,i}} \left[ \log\left(1 + \frac{\rho}{M} \cdot \|\mathbf{h}_{\mathsf{iid},\,i}\|^2 \cdot \widehat{\mathbf{h}}_{\mathsf{iid},\,i}^H \,\widetilde{\mathbf{\Lambda}}_i \,\widehat{\mathbf{h}}_{\mathsf{iid},\,i} \right) \right] \right],\tag{19}$$

where we have also used the fact that a fixed<sup>9</sup> unitary transformation of an isotropically distributed vector on  $\mathcal{G}(M, 1)$  does not alter its distribution.

## III. RATE CHARACTERIZATION: TWO-USER CASE

We now restrict attention to the special case of two-users (M = 2) and focus on computing the ergodic information-theoretic rates given in (18) and (19) in closed-form. The following theorem computes the ergodic rates as a function of the covariance matrices of the two links  $(\Sigma_1 \text{ and } \Sigma_2)$ , and the choice of beamforming vectors  $(\mathbf{w}_1 \text{ and } \mathbf{w}_2)$ .

Theorem 1: The ergodic information-theoretic rate achievable at user i (where i = 1, 2) with linear beamforming in the two-user case is given as

$$E[R_i] = E[I_{i,1}] - E[I_{i,2}] = \frac{\mathbf{\Lambda}_{i,1} h\left(\frac{\rho \mathbf{\Lambda}_{i,1}}{2}\right) - \mathbf{\Lambda}_{i,2} h\left(\frac{\rho \mathbf{\Lambda}_{i,2}}{2}\right)}{\mathbf{\Lambda}_{i,1} - \mathbf{\Lambda}_{i,2}} - h\left(\frac{\rho \widetilde{\mathbf{\Lambda}}_{i,1}}{2}\right)$$
(20)

where  $h(\bullet)$  is a monotonically increasing function defined as

$$h(x) \triangleq \exp\left(\frac{1}{x}\right) E_1\left(\frac{1}{x}\right), \ x \in (0,\infty)$$
 (21)

with  $E_1(x) = \int_x^{\infty} \frac{e^{-t}}{t} dt$  denoting the Exponential integral [44]. The corresponding eigenvalues (cf. (7) and (10)) can be written in terms of  $\Sigma_i$  and the beamforming vectors as follows:

$$\Lambda_{i,1} = \frac{A_i + B_i + \sqrt{(A_i - B_i)^2 + 4C_i^2}}{2}$$
(22)

$$\Lambda_{i,2} = \frac{A_i + B_i - \sqrt{(A_i - B_i)^2 + 4C_i^2}}{2}$$
(23)

$$\widetilde{\mathbf{\Lambda}}_{i,1} = B_i, \tag{24}$$

<sup>9</sup>Note that the unitary transformation is independent of the channel realization when the beamforming vectors are dependent only on the long-term statistics of the channel.

where  $A_i = \mathbf{w}_i^H \boldsymbol{\Sigma}_i \mathbf{w}_i$ ,  $B_i = \mathbf{w}_j^H \boldsymbol{\Sigma}_i \mathbf{w}_j$  and  $C_i = |\mathbf{w}_i^H \boldsymbol{\Sigma}_i \mathbf{w}_j|$  with  $j \neq i$  and  $\{i, j\} = 1, 2$ .

*Proof:* Since  $E[R_i] = E[I_{i,1}] - E[I_{i,2}]$ , we start by computing  $E[I_{i,1}]$ . From (18), we have

$$E\left[I_{i,1}\right] = E_{\mathbf{X}}\left[\int_{y=\mathbf{\Lambda}_{i,2}}^{\mathbf{\Lambda}_{i,1}} \log\left(1 + \frac{\rho}{2} \cdot xy\right) p_i(y) dy\right]$$
(25)

where X stands for the random variable  $\mathbf{X} = \|\mathbf{h}_{\text{iid},i}\|^2$ , x is a realization of X and  $p_i(y)$  denotes the density function of

$$\widehat{\mathbf{h}}_{\mathsf{iid},i}^{H} \mathbf{\Lambda}_{i} \widehat{\mathbf{h}}_{\mathsf{iid},i} = \sum_{j=1}^{2} \mathbf{\Lambda}_{i,j} \left| \widehat{\mathbf{h}}_{\mathsf{iid},i}(j) \right|^{2},$$
(26)

evaluated at y with  $\Lambda_{i,2} \leq y \leq \Lambda_{i,1}$ . That is, a closed-form computation of  $E[I_{i,1}]$  requires the density function of weighted-norm of vectors isotropically distributed on  $\mathcal{G}(2,1)$ . In Lemma 1 of Appendix A, we show that

$$p_i(y) = \frac{1}{\Lambda_{i,1} - \Lambda_{i,2}}, \ \Lambda_{i,2} \le y \le \Lambda_{i,1}.$$
(27)

Using this information along with the chi-squared structure of  $\|\mathbf{h}_{\text{iid},i}\|^2$  (see (15)), we have

$$E\left[I_{i,1}\right] = \frac{1}{\Lambda_{i,1} - \Lambda_{i,2}} \cdot \int_{x=0}^{\infty} x e^{-x} \int_{y=\Lambda_{i,2}}^{\Lambda_{i,1}} \log\left(1 + \frac{\rho}{2}xy\right) dy \, dx.$$
(28)

Integrating out the y variable, we have

$$E[I_{i,1}] = \frac{1}{\frac{\rho}{2} \left( \mathbf{\Lambda}_{i,1} - \mathbf{\Lambda}_{i,2} \right)} \cdot \int_{x=0}^{\infty} \left( 1 + \frac{\rho}{2} \mathbf{\Lambda}_{i,1} x \right) \cdot \log \left( 1 + \frac{\rho}{2} \mathbf{\Lambda}_{i,1} x \right) \cdot e^{-x} dx$$
$$- \frac{1}{\frac{\rho}{2} \left( \mathbf{\Lambda}_{i,1} - \mathbf{\Lambda}_{i,2} \right)} \cdot \int_{x=0}^{\infty} \left( 1 + \frac{\rho}{2} \mathbf{\Lambda}_{i,2} x \right) \cdot \log \left( 1 + \frac{\rho}{2} \mathbf{\Lambda}_{i,2} x \right) \cdot e^{-x} dx - 1.$$
(29)

Following a routine computation using the list of integral table formula [45, 4.337(2), p. 572], we have the expression for  $E[I_{i,1}]$ . Particularizing this expression to the case of  $E[I_{i,2}]$  in (19) with  $\widetilde{\Lambda}_{i,2} = 0$  results in the rate expression as in the statement of the theorem. To complete the proposition, an elementary computation of the eigenvalues of the associated  $2 \times 2$  matrices in (7) and (10) results in their characterization.

The increasing nature of  $h(\bullet)$ , defined in (21), is illustrated in Fig. 1. Towards the goal of obtaining physical intuition on the structure of the optimal beamforming vectors, it is of interest to obtain the limiting form of the ergodic rates in the low- and the high-SNR extremes.

## A. Low-SNR Extreme

Proposition 1: The ergodic rate  $E[R_i]$  can be bounded as

$$1 - \rho \mathsf{C}_{\mathsf{low}} \le \frac{E\left[R_i\right]}{\frac{\rho}{2}\left(\mathbf{\Lambda}_{i,1} + \mathbf{\Lambda}_{i,2} - B_i\right)} \le 1 + \rho \mathsf{C}_{\mathsf{up}}$$
(30)



Fig. 1. The behavior of h(x) for x satisfying  $0 < x \le 25$ .

for some positive constants  $C_{up}$  and  $C_{low}$  (not provided here for the sake of brevity) that depend only on the eigenvalues  $\Lambda_{i,1}, \Lambda_{i,2}$  and  $\widetilde{\Lambda}_{i,1}$ . Thus, as  $\rho \to 0$ , we have

$$\frac{E[R_i]}{\rho} \xrightarrow{\rho \to 0} \frac{1}{2} \left( \mathbf{\Lambda}_{i,1} + \mathbf{\Lambda}_{i,2} - B_i \right)$$
(31)

$$= \frac{A_i}{2} = \frac{\mathbf{w}_i^H \boldsymbol{\Sigma}_i \mathbf{w}_i}{2}.$$
 (32)

*Proof:* We need the following bounds on the Exponential integral [44, 5.1.20, p. 229]:

$$\frac{x}{1+2x} \le \frac{1}{2} \log (1+2x) \le h(x) \le \log (1+x) \le x$$
(33)

where the extremal inequalities are established by using the fact that

$$\frac{x}{x+1} \le \log(1+x) \le x. \tag{34}$$

Using these bounds, it is straightforward to see that the relationship in (30) holds. Note that both the upper and lower bounds converge to the same value as  $\rho \rightarrow 0$ , which results in the simplification in (32).

In the low-SNR extreme, the system is noise-limited, hence the linear scaling of  $E[R_i]$  with  $\rho$ .

#### B. High-SNR Extreme

*Proposition 2:* As  $\rho \to \infty$ , we have

$$E[R_i] \xrightarrow{\rho \to \infty} \frac{\mathbf{\Lambda}_{i,1} \log(\mathbf{\Lambda}_{i,1}) - \mathbf{\Lambda}_{i,2} \log(\mathbf{\Lambda}_{i,2})}{\mathbf{\Lambda}_{i,1} - \mathbf{\Lambda}_{i,2}} - \log(B_i)$$
(35)

$$= \frac{A_i + B_i}{2\sqrt{(A_i - B_i)^2 + 4C_i^2}} \cdot \log\left(\frac{A_i + B_i + \sqrt{(A_i - B_i)^2 + 4C_i^2}}{A_i + B_i - \sqrt{(A_i - B_i)^2 + 4C_i^2}}\right) + \frac{1}{2}\log\left(\frac{A_iB_i - C_i^2}{B_i^2}\right)$$
(36)

with  $A_i, B_i$  and  $C_i$  as in Theorem 1.

*Proof:* The following asymptotic expansion [44, 5.1.11, p. 229] of the Exponential integral is useful in obtaining the limiting form of  $E[R_i]$  as  $\rho \to \infty$ :

$$E_1(x) = \log\left(\frac{1}{x}\right) + \sum_{k=1}^{\infty} \frac{(-1)^{k+1} x^k}{k \cdot k!} - \gamma$$
(37)

$$\stackrel{x \to 0}{\to} \log\left(\frac{1}{x}\right) + x - \gamma \tag{38}$$

where  $\gamma \approx 0.577$  is the Euler-Mascheroni constant. Using the limiting value of  $E_1(x)$  to approximate  $E[R_i]$ , we have the expression in (35). Expanding  $\Lambda_{i,1}$  and  $\Lambda_{i,2}$  in terms of  $A_i, B_i$  and  $C_i$ , we have the expression in (36).

Unlike the low-SNR extreme,  $E[R_i]$  is not a function of  $\rho$  here. The dominating impact of interference (due to the fixed nature of the beamforming vectors that are not adapted to the channel realizations) and the consequent boundedness of  $E[R_i]$  in (35) as  $\rho$  increases should not be surprising.

## IV. SUM-RATE OPTIMIZATION: TWO-USER CASE

We are now interested in understanding the structure of the optimal choice of beamforming vectors  $(\mathbf{w}_{1, \text{opt}}, \mathbf{w}_{2, \text{opt}})$  that maximize  $\mathcal{R}_{sum}$  as a function of  $\Sigma_1$ ,  $\Sigma_2$  and  $\rho$ . This problem is difficult, in general. To obtain insight, we first consider the low- and the high-SNR extremes before studying the intermediate-SNR regime.

For simplicity, let us assume an eigen-decomposition for  $\Sigma_1$  and  $\Sigma_2$  of the form

$$\Sigma_1 = \mathbf{U} \operatorname{diag}([\lambda_1(\Sigma_1), \lambda_2(\Sigma_1)]) \mathbf{U}^H,$$
(39)

$$\Sigma_2 = \widetilde{\mathbf{U}} \operatorname{diag}([\lambda_1(\Sigma_2), \lambda_2(\Sigma_2)]) \widetilde{\mathbf{U}}^H,$$
(40)

where  $\mathbf{U} = [\mathbf{u}_1(\boldsymbol{\Sigma}_1), \mathbf{u}_2(\boldsymbol{\Sigma}_1)], \quad \widetilde{\mathbf{U}} = [\mathbf{u}_1(\boldsymbol{\Sigma}_2), \mathbf{u}_2(\boldsymbol{\Sigma}_2)], \text{ and } \lambda_1(\boldsymbol{\Sigma}_i) \geq \lambda_2(\boldsymbol{\Sigma}_i), \quad i = 1, 2.$  In particular, we assume that both  $\boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2$  are positive-definite, that is,  $\lambda_2(\boldsymbol{\Sigma}_i) > 0$ .

### A. Low-SNR Extreme

*Proposition 3:* In the low-SNR regime, from Prop. 1 we see that the maximization of  $E[R_i]$  involves optimizing over  $\mathbf{w}_i$  alone. Thus, we have

$$\mathbf{w}_{i, \text{opt}} = \arg\max_{\mathbf{w}_i} \mathcal{R}_{\text{sum}} = \arg\max_{\mathbf{w}_i} \mathbf{w}_i^H \boldsymbol{\Sigma}_i \mathbf{w}_i = e^{j\nu_i} \mathbf{u}_1(\boldsymbol{\Sigma}_i)$$
(41)

for some choice of  $\nu_i \in [0, 2\pi)$ , i = 1, 2. The resulting ergodic sum-rate satisfies

$$\lim_{\rho \to 0} \frac{\mathcal{R}_{\mathsf{sum}}}{\rho} = \frac{1}{2} \cdot \left[ \lambda_1(\Sigma_1) + \lambda_1(\Sigma_2) \right].$$
(42)

In the low-SNR extreme, the optimal solution is such that the transmitter signals to a given user along the dominant statistical eigen-mode of that user's channel and ignores the other user's channel completely. This is a solution motivated by the single-user viewpoint where the optimality of signalling along the dominant statistical eigen-mode of the forward channel is well-known [1]–[5]. This solution is not surprising since in the noise-limited regime, the broadcast channel is well-approximated by separate single-user models connecting the transmitter to each receiver.

#### B. High-SNR Extreme

Define  $\Sigma$  (and its corresponding eigen-decomposition) as

$$\boldsymbol{\Sigma} \triangleq \boldsymbol{\Sigma}_{2}^{-\frac{1}{2}} \boldsymbol{\Sigma}_{1} \boldsymbol{\Sigma}_{2}^{-\frac{1}{2}} = \mathbf{V} \operatorname{diag}\left( \left[ \eta_{1} \ \eta_{2} \right] \right) \mathbf{V}^{H}$$
(43)

where  $\mathbf{V} = [\mathbf{v}_1 \ \mathbf{v}_2]$  and  $\eta_1 \ge \eta_2$ . Note that  $\Sigma$  is positive-definite ( $\eta_2 > 0$ ) since both  $\Sigma_1$  and  $\Sigma_2$  are positive-definite. The main result of this section is as follows.

*Theorem 2:* In the high-SNR extreme, the ergodic sum-rate is maximized by the following choice of beamforming vectors:

$$\mathbf{w}_{1,\,\mathsf{opt}} = e^{j\nu_1} \cdot \frac{\mathbf{\Sigma}_2^{-\frac{1}{2}} \mathbf{v}_1}{\|\mathbf{\Sigma}_2^{-\frac{1}{2}} \mathbf{v}_1\|}, \quad \mathbf{w}_{2,\,\mathsf{opt}} = e^{j\nu_2} \cdot \frac{\mathbf{\Sigma}_2^{-\frac{1}{2}} \mathbf{v}_2}{\|\mathbf{\Sigma}_2^{-\frac{1}{2}} \mathbf{v}_2\|}$$
(44)

for some choice of  $\nu_i \in [0, 2\pi)$ , i = 1, 2. The optimal ergodic sum-rate satisfies

$$\lim_{\rho \to \infty} \mathcal{R}_{\mathsf{sum}} = \frac{\kappa_1 \log (\kappa_1)}{\kappa_1 - 1} + \frac{\log (\kappa_2)}{\kappa_2 - 1}$$
(45)

where

$$\kappa_1 \triangleq \frac{\eta_1 \tau_2}{\eta_2 \tau_1}, \qquad \kappa_2 \triangleq \frac{\tau_2}{\tau_1},$$
(46)

$$\tau_1 = \mathbf{v}_1^H \boldsymbol{\Sigma}_2^{-1} \mathbf{v}_1, \quad \tau_2 = \mathbf{v}_2^H \boldsymbol{\Sigma}_2^{-1} \mathbf{v}_2, \quad \tau_3 = \mathbf{v}_1^H \boldsymbol{\Sigma}_2^{-1} \mathbf{v}_2.$$
(47)

*Proof:* The first step in our proof is to rewrite the high-SNR rate expression in a form that permits further analysis. This is done in Appendices B and C. With the definition of  $\{v_1, v_2\}$  as in (43), since  $\Sigma_2$  is full-rank, we can decompose  $w_1$  and  $w_2$  as

$$\mathbf{w}_{1} = \frac{\alpha \Sigma_{2}^{-\frac{1}{2}} \mathbf{v}_{1} + \beta \Sigma_{2}^{-\frac{1}{2}} \mathbf{v}_{2}}{\|\alpha \Sigma_{2}^{-\frac{1}{2}} \mathbf{v}_{1} + \beta \Sigma_{2}^{-\frac{1}{2}} \mathbf{v}_{2}\|}$$
(48)

$$\mathbf{w}_{2} = \frac{\gamma \Sigma_{2}^{-\frac{1}{2}} \mathbf{v}_{1} + \delta \Sigma_{2}^{-\frac{1}{2}} \mathbf{v}_{2}}{\|\gamma \Sigma_{2}^{-\frac{1}{2}} \mathbf{v}_{1} + \delta \Sigma_{2}^{-\frac{1}{2}} \mathbf{v}_{2}\|}$$
(49)

for some choice of  $\{\alpha, \beta, \gamma, \delta\}$  with  $\alpha = |\alpha|e^{j\theta_{\alpha}}$  (similarly, for other quantities) satisfying  $|\alpha|^2 + |\beta|^2 = |\gamma|^2 + |\delta|^2 = 1$ . In Appendix D, we show that the ergodic sum-rate optimization over the six-dimensional parameter space  $\{|\alpha|, |\gamma|, \theta_{\alpha}, \theta_{\beta}, \theta_{\gamma}, \theta_{\delta}\}$  results in the choice as in the statement of the theorem.

Many remarks are in order at this stage.

## Remarks:

1) Recall the definition of a generalized eigenvector:

Definition 2 (Generalized eigenvector [46]): A generalized eigenvector x (with the corresponding generalized eigenvalue  $\sigma$ ) of a pair of matrices (A, B) satisfies the relationship

$$\mathbf{A}\mathbf{x} = \sigma \mathbf{B}\mathbf{x}.$$
 (50)

In the special case where B is invertible, a generalized eigenvector of the pair (A, B) is also an eigenvector of  $B^{-1}A$ . If A and B are also positive-definite, then all the generalized eigenvalues are positive. While a unit-norm generalized eigenvector (or an eigenvector) is not unique on St(M, 1), it is unique on  $\mathcal{G}(M, 1)$ .

We decompose  $\{\mathbf{w}_1, \mathbf{w}_2\}$  in (48)-(49) along the basis<sup>10</sup>  $\{\boldsymbol{\Sigma}_2^{-\frac{1}{2}}\mathbf{v}_1, \boldsymbol{\Sigma}_2^{-\frac{1}{2}}\mathbf{v}_2\}$  instead of the more routine basis  $\{\mathbf{v}_1, \mathbf{v}_2\}$ . The reason for this peculiar choice is as follows. It turns out that  $\boldsymbol{\Sigma}_2^{-\frac{1}{2}}\mathbf{v}_1$  and  $\boldsymbol{\Sigma}_2^{-\frac{1}{2}}\mathbf{v}_2$  are the dominant generalized eigenvectors (corresponding to the largest generalized eigenvalue) of the pairs  $(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2)$  and  $(\boldsymbol{\Sigma}_2, \boldsymbol{\Sigma}_1)$ , respectively. For this claim, we use (43) to note that

$$\Sigma_{2}^{-1}\Sigma_{1} = \Sigma_{2}^{-\frac{1}{2}} \left( \Sigma_{2}^{-\frac{1}{2}} \Sigma_{1} \Sigma_{2}^{-\frac{1}{2}} \right) \Sigma_{2}^{\frac{1}{2}} = \mathbf{M} \mathbf{D} \mathbf{M}^{-1}$$
(51)

$$\Sigma_{1}^{-1}\Sigma_{2} = (\Sigma_{2}^{-1}\Sigma_{1})^{-1} = \mathbf{M}\mathbf{D}^{-1}\mathbf{M}^{-1}$$
(52)

where  $\mathbf{M} = \boldsymbol{\Sigma}_2^{-\frac{1}{2}} \mathbf{V}$  and  $\mathbf{D} = \text{diag}([\eta_1 \ \eta_2])$ . This means that we can write

$$\mathbf{w}_{1,\,\mathsf{opt}} = e^{j\nu_1} \cdot \mathbf{u}_1 \left( \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1 \right) \tag{53}$$

$$\mathbf{w}_{2,\mathsf{opt}} = e^{j\nu_2} \cdot \mathbf{u}_2 \left( \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1 \right)$$
(54)

<sup>10</sup>Note that  $\Sigma_2$  is a full-rank matrix and hence, the vectors  $\Sigma_2^{-\frac{1}{2}} \mathbf{v}_1$  and  $\Sigma_2^{-\frac{1}{2}} \mathbf{v}_2$  form a non-orthogonal basis (in general), whereas  $\{\mathbf{v}_1, \mathbf{v}_2\}$  is orthonormal.

where  $\mathbf{u}_1(\bullet)$  and  $\mathbf{u}_2(\bullet)$  are the dominant and non-dominant eigenvectors, respectively. Using the generalized eigenvector structure, it is easy to see that

$$\mathbf{u}_1 \left( \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1 \right) = \mathbf{u}_2 \left( \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_2 \right)$$
(55)

$$\mathbf{u}_2 \big( \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1 \big) = \mathbf{u}_1 \big( \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_2 \big)$$
 (56)

and thus

$$\mathbf{w}_{1,\mathsf{opt}} = e^{j\nu_1} \cdot \mathbf{u}_1 \left( \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1 \right)$$
(57)

$$\mathbf{w}_{2,\,\mathsf{opt}} = e^{j\nu_2} \cdot \mathbf{u}_1 \big( \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_2 \big).$$
(58)

2) Given that the transmitter has only statistical information of the two links, a natural candidate for beamforming in the high-SNR extreme is the solution to the maximization of an appropriately-defined average SINR metric for each user. Motivated by the fact (see (3)) that the instantaneous sum-rate for the *i*-th user ( $R_i$ ) is an increasing function of  $|\mathbf{h}_i^H \mathbf{w}_i|^2$  whereas  $R_j$  (for  $j \neq i$ ) is a decreasing function of  $|\mathbf{h}_j^H \mathbf{w}_i|^2$ , we define an "average" SINR metric as follows:

$$\mathsf{SINR}_{i} \triangleq \frac{E\left[|\mathbf{h}_{i}^{H}\mathbf{w}_{i}|^{2}\right]}{E\left[|\mathbf{h}_{j}^{H}\mathbf{w}_{i}|^{2}\right]} = \frac{\mathbf{w}_{i}^{H}\boldsymbol{\Sigma}_{i}\mathbf{w}_{i}}{\mathbf{w}_{i}^{H}\boldsymbol{\Sigma}_{j}\mathbf{w}_{i}}.$$
(59)

The optimization problem of interest is to maximize  $SINR_i$  which has the generalized eigenvector structure as solution [47]:

$$\arg\max_{\mathbf{w}_i:\mathbf{w}_i^H\mathbf{w}_i=1}\mathsf{SINR}_i = e^{j\nu_i}\mathbf{u}_1\left(\boldsymbol{\Sigma}_j^{-1}\boldsymbol{\Sigma}_i\right), \ j \neq i, \ \{i,j\} = 1, 2.$$
(60)

It follows that if user *i* selfishly maximizes (its own) SINR<sub>*i*</sub> metric, then the set of such beamforming vectors maximize the ergodic sum-rate in the high-SNR regime. In this sense, the solution to the broadcast problem mirrors and generalizes the single-user setting, where the optimality of signalling along the statistical eigen-modes of the channel is well-understood [1]–[5]. Further, while optimal beamformer solutions in terms of the generalized eigenvectors are obtained in the perfect CSI case of the broadcast setting for the beamforming design problem [26], [37] and the interference channel problem [38], to the best of our knowledge, this solution in the statistical case is a first. A similar result is obtained in a related work of ours [39] on statistical beamforming vector design for the interference channel case. Since the generalized eigenvector solution has an intuitive explanation, it is of interest to obtain useful insights on the optimality of this solution in more general multi-user settings.

3) The ergodic sum-rate in (45) is increasing in  $\kappa_1$  and thus in  $\frac{\eta_1}{\eta_2}$ . We now observe that ill-conditioning of  $\Sigma_1$  is necessary and sufficient to ensure that  $\frac{\eta_1}{\eta_2}$  is large. For this, we

use standard eigenvalue inequalities for product of Hermitian matrices [47] to see that

$$\frac{\chi_1}{\chi_2} = \frac{\lambda_1(\Sigma_1) \cdot \lambda_2(\Sigma_2^{-1})}{\lambda_2(\Sigma_1) \cdot \lambda_1(\Sigma_2^{-1})} \le \frac{\eta_1}{\eta_2} \le \frac{\lambda_1(\Sigma_1) \cdot \lambda_1(\Sigma_2^{-1})}{\lambda_2(\Sigma_1) \cdot \lambda_2(\Sigma_2^{-1})} = \chi_1 \cdot \chi_2$$
(61)

where  $\chi_i = \frac{\lambda_1(\Sigma_i)}{\lambda_2(\Sigma_i)}$ , i = 1, 2. In other words, the more ill-conditioned  $\Sigma_1$  is, the larger the high-SNR statistical beamforming sum-rate asymptote is (and *vice versa*).

On the other hand, the ergodic sum-rate in (45) is not monotonic in  $\frac{\tau_1}{\tau_2}$ . Nevertheless, it can be seen that as a function of  $\frac{\tau_1}{\tau_2}$ , it has local maxima as  $\frac{\tau_1}{\tau_2} \to 0$  and  $\frac{\tau_1}{\tau_2} \to \infty$ , and a minimum at  $\frac{\tau_1}{\tau_2} = 1$ . The more well-conditioned  $\Sigma_2$  is, the more closer  $\frac{\tau_1}{\tau_2}$  is to 1 and hence, the high-SNR statistical beamforming sum-rate asymptote is minimized. If  $\Sigma_2$  is ill-conditioned, the value taken by  $\frac{\tau_1}{\tau_2}$  depends on the angle between the dominant eigenvectors of  $\Sigma_2$  and  $\Sigma$ . If the two eigenvectors are nearly parallel,  $\frac{\tau_1}{\tau_2}$  is close to zero and if they are nearly perpendicular,  $\frac{\tau_1}{\tau_2}$  is very large. In either case, the high-SNR statistical beamforming sum-rate asymptote is locally maximized.

The conclusion from the above analysis is that among all possible channels, the ergodic sum-rate is maximized (or minimized) when  $\Sigma_1$  and  $\Sigma_2$  are both ill- (or well-)conditioned. In other words, if both the users encounter poor scattering (that leads to an ill-conditioning of their respective covariance matrices), their fading is spatially localized. The transmitter can simultaneously excite these spatial localizations without causing a proportional increase in the interference level of the other user thus resulting in a higher ergodic sum-rate. On the other hand, rich scattering implies that fading is spatially isotropic for both the users. Any spatially localized excitation for one user will cause an isotropic interference level at the other user thus resulting in a smaller ergodic sum-rate.

4) A special case that is of considerable interest is when  $\Sigma_1$  and  $\Sigma_2$  have the same set of orthonormal eigenvectors. This would be a suitable model for certain indoor scenarios where the antenna separation for the two users is the same [28]. Denoting (for simplicity) the set of common eigenvectors by  $\mathbf{u}_1$  and  $\mathbf{u}_2$ , we can decompose  $\Sigma_1$  and  $\Sigma_2$  as

$$\boldsymbol{\Sigma}_{1} = \left[\mathbf{u}_{1}, \mathbf{u}_{2}\right] \operatorname{diag}\left(\left[\lambda_{1}, \lambda_{2}\right]\right) \left[\mathbf{u}_{1}, \mathbf{u}_{2}\right]^{H}, \tag{62}$$

$$\boldsymbol{\Sigma}_{2} = \left[ \mathbf{u}_{1}, \mathbf{u}_{2} \right] \operatorname{diag} \left( \left[ \mu_{1}, \ \mu_{2} \right] \right) \left[ \mathbf{u}_{1}, \mathbf{u}_{2} \right]^{H}.$$
(63)

We re-use the notations  $\chi_1$  and  $\chi_2$  to denote

$$\chi_1 \triangleq \frac{\lambda_1}{\lambda_2} \quad \text{and} \quad \chi_2 \triangleq \frac{\mu_1}{\mu_2}.$$
 (64)

Without loss in generality, we can assume that  $\chi_1 \ge 1$ . Two scenarios<sup>11</sup> arise depending on the relationship between  $\chi_1$  and  $\chi_2$ : i)  $\chi_1 \ge \chi_2$ , and ii)  $\chi_1 < \chi_2$ .

<sup>&</sup>lt;sup>11</sup>These possibilities arise because even though the set of eigenvectors of  $\Sigma_1$  and  $\Sigma_2$  are the same, there is no specific reason to expect the dominant eigenvector of  $\Sigma_1$  to also be a dominant eigenvector of  $\Sigma_2$ . Observe that the first case subsumes the setting where  $\mu_1 = \mu_2 = \mu$  and  $\Sigma_2 = \mu I$ .

*Theorem 3:* In the high-SNR extreme, the ergodic sum-rate is maximized by the following choice of beamforming vectors:

$$\mathbf{w}_{1,\,\mathsf{opt}} = e^{j\nu_1} \mathbf{u}_1, \quad \mathbf{w}_{2,\,\mathsf{opt}} = e^{j\nu_2} \mathbf{u}_2 \quad \text{if} \quad \chi_1 \ge \chi_2, \\ \mathbf{w}_{1,\,\mathsf{opt}} = e^{j\nu_2} \mathbf{u}_2, \quad \mathbf{w}_{2,\,\mathsf{opt}} = e^{j\nu_1} \mathbf{u}_1 \quad \text{if} \quad \chi_1 < \chi_2$$
(65)

for some choice of  $\nu_i \in [0, 2\pi)$ , i = 1, 2. The optimal ergodic sum-rate satisfies

$$\lim_{\rho \to \infty} \mathcal{R}_{\mathsf{sum}} = \begin{cases} \frac{\chi_1 \cdot \log(\chi_1)}{\chi_1 - 1} + \frac{\log(\chi_2)}{\chi_2 - 1} & \text{if } \chi_1 \ge \chi_2 \\ \frac{\chi_2 \cdot \log(\chi_2)}{\chi_2 - 1} + \frac{\log(\chi_1)}{\chi_1 - 1} & \text{if } \chi_1 < \chi_2. \end{cases}$$
(66)

*Proof:* While Theorem 2 can be particularized to this special case easily, we pursue an alternate proof technique in Appendix E that exploits the comparative relationship between  $\tau_1$  and  $\tau_2$  (which is possible in the special case) and the fact that  $\tau_3 = 0$ .

- A comparison of the proof techniques of Theorems 2 and 3 is presented in Appendix F. 5) Some remarks on the optimization set-up of this paper are necessary. The proofs of Theorems 2 and 3 require us to consider a six-dimensional optimization over the parameter space of  $\{|\alpha|, |\gamma|, \theta_{\alpha}, \theta_{\beta}, \theta_{\gamma}, \theta_{\delta}\}$ . As a result, any geometric interpretation of the optimization is impossible. A naive approach to the six-dimensional optimization problems in this paper is to adopt the method of matrix differentiation calculus. For this approach to work, we need to show that both the set over which optimization is done as well as the optimized function are convex, neither of which is true in our case. Specifically, neither  $\mathcal{G}(2, 1)$  nor St(2, 1) are convex sets. It also turns out that the ergodic sum-rate is neither convex nor concave<sup>12</sup> over the set of beamforming vectors, even over a locally convex domain or an extended convex domain (like the interior of the sphere).
- 6) The approach adopted in Appendices D and E overcomes these difficulties, and it consists of two steps. In the first step, we produce an upper bound to the ergodic sum-rate that is independent of the optimization parameters. In the second step, we show that this upper bound can be realized by a specific choice of beamforming vectors thereby confirming that choice's optimality. This approach seems to be the most natural (and first principles-based) recourse to solving the non-convex optimization problem at hand. An alternate approach to optimize the ergodic sum-rate is non-linear optimization theory [48]. But this approach is fraught with complicated Hessian calculations and technical difficulties such as distinguishing between local and global extrema.

<sup>&</sup>lt;sup>12</sup>It is possible that some function of the ergodic sum-rate may be convex. But we are not aware of any likely candidate that could work.

## C. Intermediate-SNR Regime: Candidate Beamforming Vectors

While physical intuition on the structure of the optimal ergodic sum-rate maximizing beamforming vectors has been obtained in the low- and the high-SNR extremes, the essentially intractable nature of the Exponential integral in the ergodic rate expressions of Theorem 1 means that such a possibility at an arbitrary SNR is difficult. Nevertheless, the single-user set-up [49], [50] suggests that the optimal beamforming vectors (that determine the modes that are excited) and the power allocation across these modes can be continuously parameterized by a function of the SNR. Motivated by the single-user case, a desirable quality for a "good" beamforming vector structure ({ $w_{i, cand}(\rho), i = 1, 2$ }) at an arbitrary SNR of  $\rho$  is that the limiting behavior of such a structure in the low- and the high-SNR extremes should be the solutions of Prop. 3 and Theorem 2. That is,

$$\lim_{\rho \to 0} \mathbf{w}_{i, \text{ cand}}(\rho) = e^{j\nu_i} \mathbf{u}_1(\boldsymbol{\Sigma}_i), \tag{67}$$

$$\lim_{\rho \to \infty} \mathbf{w}_{i, \operatorname{cand}}(\rho) = e^{j\nu_i} \frac{\boldsymbol{\Sigma}_2^{-\frac{1}{2}} \mathbf{v}_i}{\|\boldsymbol{\Sigma}_2^{-\frac{1}{2}} \mathbf{v}_i\|}, \quad i = 1, 2,$$
(68)

where the above limits are seen as manifold operations [42] on  $\mathcal{G}(2,1)$ .

A natural candidate that meets (67) and (68) is the following choice parameterized by  $\alpha(\rho)$ and  $\beta(\rho)$  satisfying  $\{\alpha(\rho), \beta(\rho)\} \in [0, \infty)$  and  $\nu_i \in [0, 2\pi), i = 1, 2$ :

$$\mathbf{w}_{1,\,\mathsf{cand}}(\rho) = e^{j\nu_1} \cdot \mathsf{Dom.\,eig.}\left(\left(\alpha(\rho)\boldsymbol{\Sigma}_2 + \mathbf{I}\right)^{-1}\boldsymbol{\Sigma}_1\right)$$
(69)

$$\mathbf{w}_{2,\,\mathsf{cand}}(\rho) = e^{j\nu_2} \cdot \mathsf{Dom.\,eig.}\left(\left(\beta(\rho)\boldsymbol{\Sigma}_1 + \mathbf{I}\right)^{-1}\boldsymbol{\Sigma}_2\right)$$
(70)

where the notation Dom.  $eig(\bullet)$  stands for the unit-norm dominant eigenvector operation. These vectors can be seen to be solutions to the following optimization problems:

$$\arg\max_{\mathbf{w}_i: \mathbf{w}^H \mathbf{w}_i = 1} \mathsf{SINR}_i = \mathbf{w}_{i, \mathsf{cand}}(\rho), \ i = 1, 2$$
(71)

where

$$SINR_{1} = \frac{\mathbf{w}_{1}^{H} \boldsymbol{\Sigma}_{1} \mathbf{w}_{1}}{\mathbf{w}_{1}^{H} \mathbf{w}_{1} + \alpha(\rho) \mathbf{w}_{1}^{H} \boldsymbol{\Sigma}_{2} \mathbf{w}_{1}}$$
(72)

$$SINR_2 = \frac{\mathbf{w}_2^H \boldsymbol{\Sigma}_2 \mathbf{w}_2}{\mathbf{w}_2^H \mathbf{w}_2 + \beta(\rho) \mathbf{w}_2^H \boldsymbol{\Sigma}_1 \mathbf{w}_2}.$$
(73)

The choice in (69)-(70) is a low-dimensional mapping from  $\mathcal{G}(2,1) \times \mathcal{G}(2,1)$  to  $\mathbb{R}^+ \times \mathbb{R}^+$  thus considerably simplifying the search space for candidate beamforming vectors. It must be noted that while the search space is simplified, the generalized eigenvector operation is a non-linear mapping [46] in  $\alpha(\rho)$  and  $\beta(\rho)$ .

тт

#### D. Numerical Studies

We now study the ergodic sum-rate performance with  $w_1$  and  $w_2$  as in (69)-(70) via two numerical examples. In the first study, we consider a system (note that  $Tr(\Sigma_1) = Tr(\Sigma_2) = M = 2$ ) with

$$\Sigma_{1} = \begin{bmatrix} 1.7745 & -0.5178 + 0.0247i \\ -0.5178 - 0.0247i & 0.2255 \end{bmatrix},$$
(74)  
$$\Sigma_{2} = \begin{bmatrix} 1.2522 & -0.8739 - 0.2711i \\ -0.8739 + 0.2711i & 0.7478 \end{bmatrix}.$$
(75)

Fig. 2(a) shows the ergodic sum-rate as a function of  $\rho$  for four schemes. For the first scheme, for every  $\rho$ , an optimal choice  $\{\alpha^*(\rho), \beta^*(\rho)\}$  is obtained from the search space  $\alpha(\rho) \times \beta(\rho) \in [0, \infty) \times [0, \infty)$  as follows:

$$\{\alpha^{\star}(\rho), \ \beta^{\star}(\rho)\} = \arg \max_{\{\alpha(\rho), \ \beta(\rho)\}} E[R_1] + E[R_2]$$
  
s.t.  $\mathbf{w}_1 = \mathbf{w}_{1, \text{ cand}}(\rho), \ \mathbf{w}_2 = \mathbf{w}_{2, \text{ cand}}(\rho).$  (76)

The performance of the beamforming vectors with  $\alpha^*(\rho)$  and  $\beta^*(\rho)$  for every  $\rho$  is plotted along with the performance of the candidate obtained via a numerical (Monte Carlo) search over  $\mathcal{G}(2,1) \times \mathcal{G}(2,1)$ . As motivated in the prior discussion, while we expect the performance with  $\{\alpha^*(\rho), \beta^*(\rho)\}$  to be good, it is surprising that this choice is indeed optimal. Further, the performance of a set of beamforming vectors with  $\alpha(\rho) = \beta(\rho) = 0$  and  $\alpha(\rho) = 100, \beta(\rho) = 15$ (fixed for all  $\rho$  in (69)-(70)) are also plotted. Observe that these two sets approximate the lowand the high-SNR solutions of Prop. 3 and Theorem 2, respectively.

In the second study, we consider a system (again, note that  $\mathsf{Tr}(\Sigma_1) = \mathsf{Tr}(\Sigma_2) = M = 2$ ) with

$$\Sigma_{1} = \begin{bmatrix} 1.3042 & 0.0543 - 0.2540i \\ 0.0543 + 0.2540i & 0.6958 \end{bmatrix},$$
(77)

$$\Sigma_2 = \begin{bmatrix} 1.1161 & -0.2195 + 0.4340i \\ -0.2195 - 0.4340i & 0.8839 \end{bmatrix}.$$
 (78)

Fig. 2(b) plots the performance of the proposed scheme, the low- and the high-SNR solutions in addition to the candidate obtained via a numerical search over  $\mathcal{G}(2,1) \times \mathcal{G}(2,1)$ . As before, the low- and the high-SNR solutions are optimal in their respective extremes while the candidate  $\{\alpha^*(\rho), \beta^*(\rho)\}$  is optimal across all  $\rho$ .

Note that in Fig. 2(a) there exists an SNR-regime where both the low- and the high-SNR solutions are sub-optimal. In contrast, in Fig. 2(b), the high-SNR solution essentially coincides with the numerical search for all  $\rho$  whereas at the low-SNR extreme, the performance of the



Fig. 2. Performance of proposed scheme with  $\Sigma_1$  and  $\Sigma_2$ : (a) as in (74)-(75), (b) as in (77)-(78).

low-SNR solution is as expected. We now explain why the high-SNR solution performs as well as  $\{\alpha^*(\rho), \beta^*(\rho)\}\$  for all  $\rho$ . For this, we need to understand the behavior of the angle between the proposed set of beamforming vectors in (69)-(70) and the low-SNR solution as a function of  $\alpha$  and  $\beta$ . In Fig. 3(b), we plot  $\cos(\text{Angle}_1(\alpha(\rho)))$  as a function of  $\alpha(\rho)$  and  $\cos(\text{Angle}_2(\beta(\rho)))$ as a function of  $\beta(\rho)$  where

$$\cos\left(\mathsf{Angle}_{1}(\alpha(\rho))\right) = \left| \left(\mathsf{Dom.\,eig.}\left(\left(\alpha(\rho)\boldsymbol{\Sigma}_{2} + \mathbf{I}\right)^{-1}\boldsymbol{\Sigma}_{1}\right)\right)^{H}\mathsf{Dom.\,eig.}\left(\boldsymbol{\Sigma}_{1}\right) \right|$$
(79)

$$\cos\left(\mathsf{Angle}_{2}(\beta(\rho))\right) = \left| \left(\mathsf{Dom.\,eig.}\left(\left(\beta(\rho)\boldsymbol{\Sigma}_{1} + \mathbf{I}\right)^{-1}\boldsymbol{\Sigma}_{2}\right)\right)^{H}\mathsf{Dom.\,eig.}\left(\boldsymbol{\Sigma}_{2}\right) \right|.$$
(80)

From Fig. 3(b), we note that the chordal distance<sup>13</sup> between the low- and the high-SNR solutions is small (on the order of 0.05). Also, observe that there is a quick convergence of (69)-(70) as  $\alpha$  (or  $\beta$ ) increases to the high-SNR solution and hence the high-SNR solution is a good approximation to the choice { $\alpha^*(\rho), \beta^*(\rho)$ } over a large SNR range. On the other hand, from Fig. 3(a), we see that the chordal distance between the low- and the high-SNR solutions is large (on the order of 0.90), which translates to the sub-optimality gap in Fig. 2(a).

While we are unable to prove the optimality structure of the proposed scheme in the intermediate-SNR regime, motivated by our numerical studies, we make the following conjecture.

Conjecture 1: In the intermediate-SNR regime, the ergodic sum-rate is maximized by the

 $<sup>^{13}</sup>$ In short, the chordal distance is the square-root of the difference of 1 and the square of the quantity computed in (79) (or (80)). See (136) for more details.



Fig. 3. Angle<sub>1</sub> and Angle<sub>2</sub>, defined in (79)-(80), as a function of  $\alpha$  and  $\beta$  for the setting in: (a) (74)-(75), (b) (77)-(78).

following choice of beamforming vectors:

$$\mathbf{w}_{1,\mathsf{opt}} = e^{j\nu_1} \cdot \mathsf{Dom.\,eig.}\left(\left(\alpha^*(\rho)\boldsymbol{\Sigma}_2 + \mathbf{I}\right)^{-1}\boldsymbol{\Sigma}_1\right)$$
(81)

$$\mathbf{w}_{2, \text{opt}} = e^{j\nu_2} \cdot \text{Dom. eig.} \left( \left( \beta^{\star}(\rho) \boldsymbol{\Sigma}_1 + \mathbf{I} \right)^{-1} \boldsymbol{\Sigma}_2 \right)$$
(82)

for some choice of  $\nu_i \in [0, 2\pi)$ , i = 1, 2. The notation Dom.eig(•) stands for the unit-norm dominant eigenvector operation and

$$\{\alpha^{\star}(\rho), \ \beta^{\star}(\rho)\} = \arg \max_{\{\alpha(\rho), \ \beta(\rho)\}} E[R_1] + E[R_2]$$
  
s.t.  $\mathbf{w}_1 = \mathbf{w}_{1, \text{cand}}(\rho), \ \mathbf{w}_2 = \mathbf{w}_{2, \text{cand}}(\rho).$  (83)

## V. ERGODIC SUM-RATE: GENERALIZATIONS

We studied the structure of ergodic sum-rate maximizing beamforming vectors in Sec. IV. In this section, we consider more general problems of this nature.

## A. Maximizing $E[R_i]$

Consider a system where the Quality-of-Service metric of one user significantly dominates that of the other user. For example, one user is considerably more important to the network operator than the other. The relevant metric to optimize in this scenario is not the ergodic sum-rate, but the rate achievable by the more important user. In this setting, we have the following result.

*Proposition 4:* The optimal choice of the pair  $(\mathbf{w}_{1, opt}, \mathbf{w}_{2, opt})$  that maximizes  $E[R_i]$  is: *i)* Low-SNR Extreme:

$$\mathbf{w}_{i, \mathsf{opt}} = e^{j\nu_1} \mathbf{u}_1(\mathbf{\Sigma}_i) \text{ and } \mathbf{w}_{j, \mathsf{opt}} = \mathsf{any vector on } \mathcal{G}(2, 1), \ j \neq i$$
 (84)

ii) High-SNR Extreme:

$$\mathbf{w}_{i,\,\mathsf{opt}} = e^{j\nu_1} \mathbf{u}_1\left(\mathbf{\Sigma}_i\right) \quad \text{and} \quad \mathbf{w}_{j,\,\mathsf{opt}} = e^{j\nu_2} \mathbf{u}_2\left(\mathbf{\Sigma}_i\right), \ j \neq i$$
(85)

for some choice of  $\nu_i \in [0, 2\pi), i = 1, 2$ .

*Proof:* See Appendix G.

With the above choice of beamforming vectors,  $E[R_i]$  can be written as

$$E[R_i] \stackrel{\rho \to \infty}{\to} \frac{\chi_i \log(\chi_i)}{\chi_i - 1}$$
(86)

where  $\chi_i = \frac{\lambda_1(\Sigma_i)}{\lambda_2(\Sigma_i)}$ . From (86), it is to be noted that  $E[R_i]$  increases as  $\chi_i$  increases. That is, the more ill-conditioned  $\Sigma_i$  is, the larger the high-SNR statistical beamforming rate asymptote is (and *vice versa*). This should be intuitive as our goal is only to maximize  $E[R_i]$  and the beamforming vectors in (85) achieve that goal.

### B. Weighted Ergodic Sum-Rate Maximization

In a system where the Quality-of-Service metrics of the two users are comparable (but not the same), the relevant metric to optimize is the weighted-sum of ergodic rates achievable by the two users [40]. Specifically, the objective function here is

$$\mathcal{R}_{\text{weighted}} = \zeta_1 E\left[R_1\right] + \zeta_2 E\left[R_2\right] \tag{87}$$

for some choice of weights  $\zeta_1$  and  $\zeta_2$  satisfying (without loss in generality)  $\{\zeta_1, \zeta_2\} \in [0, 1]$ . Note that  $E[R_i]$  is a special case of this objective function with  $\zeta_1 = 1, \zeta_2 = 0$  or  $\zeta_1 = 0, \zeta_2 = 1$ .

Maximizing  $\mathcal{R}_{\text{weighted}}$  to obtain a closed-form characterization of the optimal beamforming vectors seems hard in general. Motivated by the study for the sum-rate in the intermediate-SNR regime in Sec. IV, we now consider a set of candidate beamforming vectors that produce known optimal structures in special cases. For this, it is important to note that no choice of  $\alpha(\rho)$  and  $\beta(\rho)$  in (69)-(70) can produce the beamforming vectors in (85). A candidate set of beamforming vectors that not only produces the special (extreme) cases in the ergodic sum-rate setting, but also (85) is the following choice parameterized by four quantities,  $\{\alpha(\rho), \beta(\rho), \gamma(\rho), \delta(\rho)\} \in [0, \infty)$ :

$$\mathbf{w}_{1,\text{weighted, cand}}(\rho) = e^{j\nu_1} \cdot \text{Dom. eig.} \left( \left( \alpha(\rho) \boldsymbol{\Sigma}_2 + \mathbf{I} \right)^{-1} \left( \gamma(\rho) \boldsymbol{\Sigma}_1 + \mathbf{I} \right) \right)$$
(88)

$$\mathbf{w}_{2, \text{ weighted, cand}}(\rho) = e^{j\nu_2} \cdot \text{Dom. eig.}\left(\left(\beta(\rho)\boldsymbol{\Sigma}_1 + \mathbf{I}\right)^{-1}\left(\delta(\rho)\boldsymbol{\Sigma}_2 + \mathbf{I}\right)\right)$$
(89)

where the notation Dom.eig(•) stands for the usual unit-norm dominant eigenvector operation. As before, (88)-(89) corresponds to a low-dimensional map from  $\{\mathcal{G}(2,1)\}^4$  to  $\{[0,\infty)\}^4$  and thus a simplification in the search for a good beamformer structure.



Fig. 4. Weighted ergodic sum-rate of proposed scheme with  $\Sigma_1$  and  $\Sigma_2$  as in (74)-(75).

We now study the performance of the proposed beamforming vectors in (88)-(89) for the system with  $\Sigma_1$  and  $\Sigma_2$  as in (74)-(75). Two sets of weights are considered: i)  $\zeta_1 = 1, \zeta_2 = 0.5$  and ii)  $\zeta_1 = 0.2, \zeta_2 = 0.8$ . Fig. 4 plots the performance of two schemes. The first scheme corresponds to a Monte Carlo search over  $\mathcal{G}(2,1)$ , whereas the second scheme corresponds to the use of an optimal choice  $\{\alpha^*(\rho), \beta^*(\rho), \gamma^*(\rho), \delta^*(\rho)\}$  (for every  $\rho$ ) with

$$\{\alpha^{\star}(\rho), \ \beta^{\star}(\rho), \ \gamma^{\star}(\rho), \ \delta^{\star}(\rho)\} = \arg \max_{\{\alpha(\rho), \ \beta(\rho), \ \gamma(\rho), \ \delta(\rho)\}} E[R_1] + E[R_2]$$
  
s.t.  $\mathbf{w}_1 = \mathbf{w}_{1, \text{ weighted, cand}}(\rho), \ \mathbf{w}_2 = \mathbf{w}_{2, \text{ weighted, cand}}(\rho).$  (90)

As can be seen from Fig. 4, the proposed scheme in (88)-(89) performs as well as the Monte Carlo search for both sets of weights. Numerical studies suggest that similar performance is seen across all possible  $\Sigma_1$  and  $\Sigma_2$ , and all possible weights  $\zeta_1$  and  $\zeta_2$ . Motivated by these studies, we pose the following conjecture.

Conjecture 2: In the intermediate-SNR regime, the weighted ergodic sum-rate,  $\mathcal{R}_{weighted}$ , is

maximized by the following choice of beamforming vectors:

$$\mathbf{w}_{1} = e^{j\nu_{1}} \cdot \mathsf{Dom.\,eig.}\left(\left(\alpha^{\star}(\rho)\boldsymbol{\Sigma}_{2} + \mathbf{I}\right)^{-1}\left(\gamma^{\star}(\rho)\boldsymbol{\Sigma}_{1} + \mathbf{I}\right)\right)$$
(91)

$$\mathbf{w}_{2} = e^{j\nu_{2}} \cdot \text{Dom.eig.}\left(\left(\beta^{\star}(\rho)\boldsymbol{\Sigma}_{1} + \mathbf{I}\right)^{-1}\left(\delta^{\star}(\rho)\boldsymbol{\Sigma}_{2} + \mathbf{I}\right)\right)$$
(92)

for some choice of  $\nu_i \in [0, 2\pi)$ , i = 1, 2 and where  $\{\alpha^*(\rho), \beta^*(\rho), \gamma^*(\rho), \delta^*(\rho)\}$  are as in (90).

#### C. Rank-Deficient Case

Following up on Remark 3 in Sec. IV-B, we now consider the extreme case where both<sup>14</sup>  $\Sigma_1$  and  $\Sigma_2$  are rank-deficient in more detail.

Proposition 5: The ergodic information-theoretic rate achievable at user i is

$$E[R_i] = h\left(\frac{\rho}{2}\lambda_1(\boldsymbol{\Sigma}_i)\left(|\mathbf{u}_1(\boldsymbol{\Sigma}_i)^H\mathbf{w}_i|^2 + |\mathbf{u}_1(\boldsymbol{\Sigma}_i)^H\mathbf{w}_j|^2\right)\right) - h\left(\frac{\rho}{2}\lambda_1(\boldsymbol{\Sigma}_i)|\mathbf{u}_1(\boldsymbol{\Sigma}_i)^H\mathbf{w}_j|^2\right),$$
  
$$j \neq i, \ i = 1, 2$$
(93)

where  $h(\bullet)$  is as in (21) and the eigen-decomposition of  $\Sigma_i$  is

$$\boldsymbol{\Sigma}_{i} = \lambda_{1} (\boldsymbol{\Sigma}_{i}) \cdot \mathbf{u}_{1} (\boldsymbol{\Sigma}_{i}) \mathbf{u}_{1} (\boldsymbol{\Sigma}_{i})^{H}, \ i = 1, 2.$$
(94)

*Proof:* While Theorem 1 (as stated) is explicitly dependent on both  $\Sigma_1$  and  $\Sigma_2$  being of full rank and is hence not directly applicable in this extreme setting, much of the analysis follows through. The key to the proof is that all the results in Appendix A (Lemmas 1 and 2) also hold when some of the diagonal entries of  $\Lambda_i$  are zero. In fact, this fact is implicitly used to compute  $E[I_{i,2}]$  in Theorem 1.

## D. Three-User Case: M = 3

We now consider the task of generalizing Theorem 1 to the three-user (M = 3) case.

Proposition 6: The ergodic information-theoretic rate achievable at user i (where i = 1, 2, 3) with linear beamforming in the three-user case is

$$E[R_{i}] = E[I_{i,1}] - E[I_{i,2}]$$

$$= \frac{\Lambda_{i,1}^{2} \cdot h\left(\frac{\rho\Lambda_{i,1}}{3}\right)}{(\Lambda_{i,1} - \Lambda_{i,2})(\Lambda_{i,1} - \Lambda_{i,3})} - \frac{\Lambda_{i,2}^{2} \cdot h\left(\frac{\rho\Lambda_{i,2}}{3}\right)}{(\Lambda_{i,1} - \Lambda_{i,2})(\Lambda_{i,2} - \Lambda_{i,3})} + \frac{\Lambda_{i,3}^{2} \cdot h\left(\frac{\rho\Lambda_{i,3}}{3}\right)}{(\Lambda_{i,1} - \Lambda_{i,3})(\Lambda_{i,2} - \Lambda_{i,3})}$$

$$+ \frac{\tilde{\Lambda}_{i,1}^{2} \cdot h\left(\frac{\rho\tilde{\Lambda}_{i,1}}{3}\right)}{\left(\tilde{\Lambda}_{i,1} - \tilde{\Lambda}_{i,2}\right)\left(\tilde{\Lambda}_{i,1} - \tilde{\Lambda}_{i,2}\right)} - \frac{\tilde{\Lambda}_{i,2}^{2} \cdot h\left(\frac{\rho\tilde{\Lambda}_{i,2}}{3}\right)}{\left(\tilde{\Lambda}_{i,1} - \tilde{\Lambda}_{i,3}\right)} + \frac{\tilde{\Lambda}_{i,3}^{2} \cdot h\left(\frac{\rho\tilde{\Lambda}_{i,3}}{3}\right)}{\left(\tilde{\Lambda}_{i,1} - \tilde{\Lambda}_{i,3}\right)\left(\tilde{\Lambda}_{i,2} - \tilde{\Lambda}_{i,3}\right)}$$

$$(95)$$

where  $h(\bullet)$  is as in (21). The eigenvalue matrices  $\Lambda_i = \text{diag}([\Lambda_{i,1}, \Lambda_{i,2}, \Lambda_{i,3}])$  and  $\widetilde{\Lambda}_i = \text{diag}([\widetilde{\Lambda}_{i,1}, \widetilde{\Lambda}_{i,2}, \widetilde{\Lambda}_{i,3}])$  are defined as in (7) and (10), and can be obtained in terms of the beamforming vectors and the covariance matrices by solving the associated cubic equations.

*Proof:* The proof is tedious, but follows along the lines of Theorem 1. The first step is in characterizing  $p_i(y)$ , which is done in Lemma 2 of Appendix A. We can then generalize (28) using (121) as

$$E[I_{i,1}] = \frac{I_1}{(\Lambda_{i,1} - \Lambda_{i,2})(\Lambda_{i,1} - \Lambda_{i,3})} + \frac{I_2}{(\Lambda_{i,1} - \Lambda_{i,3})(\Lambda_{i,2} - \Lambda_{i,3})}$$
(96)

$$I_{1} = \int_{x=0}^{\infty} x^{2} e^{-x} \int_{y=0}^{\Lambda_{i,1}-\Lambda_{i,2}} y \log\left(1+\frac{\rho}{3}\Lambda_{i,1}x-\frac{\rho}{3}xy\right) dy dx$$
(97)

$$I_{2} = \int_{x=0}^{\infty} x^{2} e^{-x} \int_{y=0}^{\Lambda_{i,2}-\Lambda_{i,3}} y \log\left(1+\frac{\rho}{3}\Lambda_{i,3}x+\frac{\rho}{3}xy\right) dy \, dx.$$
(98)

These integrals are cumbersome, but straightforward to compute using [45, 4.337(2), 4.337(5), p. 572]. The result is the expression in the statement of the proposition.

#### E. General M-User Case

As can be seen from Appendix A (Lemma 2), the expression for  $p_i(y)$  becomes more complicated as M increases. Without a recourse to  $p_i(y)$ , closed-form expressions for the ergodic sumrate of the linear beamforming scheme can be obtained using a recent advance in [31], [32] that allows the computation of the density function of weighted-sum of standard central chi-squared terms (*generalized* chi-squared random variables). For example, if  $\Lambda_i = \text{diag}([\Lambda_{i,1}, \dots, \Lambda_{i,M}])$ and  $\Lambda_{i,j}$ ,  $j = 1, \dots, M$  are distinct<sup>15</sup>, we have

$$E\left[I_{i,1}\right] = \sum_{k=1}^{M} \prod_{j=1, \, j \neq k}^{M} \frac{\mathbf{\Lambda}_{i,k}}{\mathbf{\Lambda}_{i,k} - \mathbf{\Lambda}_{i,j}} \cdot h\left(\frac{\rho \mathbf{\Lambda}_{i,k}}{M}\right).$$
(99)

For  $E[I_{i,2}]$ , we replace  $\{\Lambda_{i,k}\}$  with  $\{\widetilde{\Lambda}_{i,k}\}$ . It can be checked that these expressions match with the expressions in this paper for the M = 2 and M = 3 settings. Nevertheless, it is important to note that the formulas in (95) and (99) are in terms of the eigenvalue matrices  $\{\Lambda_i, \widetilde{\Lambda}_i, i = 1, \dots, M\}$ , which become harder (and impossible for  $M \ge 5$ ) to compute in closed-form as a function of the beamforming vectors and the transmit covariance matrices as Mincreases. Tractable approximations to the ergodic sum-rate and beamforming vector optimization based on such approximations are necessary, which is the subject of ongoing work.

<sup>&</sup>lt;sup>15</sup>More complicated expressions can be obtained in case  $\{\Lambda_{i,j}\}$  are not distinct. These expressions can be derived in a straightforward manner using the results in [31].

## VI. CONCLUDING REMARKS

This paper considered the design of statistical beamforming vectors in a MISO broadcast setting to maximize the ergodic sum-rate. The approach pursued here for the simplest non-trivial problem with two-users is as follows: first, the beamforming vectors are fixed and ergodic rate expressions are computed in closed-form in terms of the covariance matrices of the links and the beamforming vectors. The optimization of this non-convex function results in a general-ized eigenvector structure for the optimal beamforming vectors, the solution to maximizing an appropriately-defined SINR metric for each user. This structure generalizes the single-user setup where the dominant eigen-modes of the transmit covariance matrix of the links are excited. The main results of this paper are presented in Table I for different SNR ( $\rho$ ) assumptions where we use  $\mathbf{u}_1(\bullet)$  and  $\mathbf{u}_2(\bullet)$  to denote the dominant and sub-dominant eigenvectors of the matrix under consideration.

Possible extensions of this work include unifying the special case of Theorem 3 with the general case of Theorem 2, and proving Conjectures 1 and 2. Developing intuition in the threeuser case as well as tractable approximations in the general M-user (M > 2) case critically depend on exploiting the functional structure of the ergodic sum-rate expression. The generalized eigenvector solution has been seen in other multi-user scenarios as well, for example, the interference channel problem with two antennas [39]. Generalizing the theme developed in the broadcast setting to the interference channel setting, the Rayleigh case to the Ricean case, and the perfect CSI case to the case where only statistical information is available are all important tasks.

Table I: Structure of Optimal Beamforming Vectors		
Objective Function:	$\arg\max_{\mathbf{w}_1,\mathbf{w}_2} E\left[R_i\right], \ i = 1, 2$	$\arg\max_{\mathbf{w}_1,\mathbf{w}_2} E\left[R_1\right] + E\left[R_2\right]$
$\rho \rightarrow 0$	$\mathbf{w}_{i,opt} = \mathbf{u}_1\left(\mathbf{\Sigma}_i ight)$	$\mathbf{w}_{1,opt} = \mathbf{u}_1\left(\mathbf{\Sigma}_1 ight)$
	$\mathbf{w}_{j, opt} = any \text{ vector on } \mathcal{G}(2, 1),$	$\mathbf{w}_{2,opt} = \mathbf{u}_1\left(\mathbf{\Sigma}_2 ight)$
	$j \neq i$ (See Prop. 4)	(See Prop. 3)
$\rho$ intermediate		$\mathbf{w}_{1,opt} = \mathbf{u}_1 \left( \left( lpha^\star( ho) \mathbf{\Sigma}_2 + \mathbf{I}  ight)^{-1} \mathbf{\Sigma}_1  ight)$
	_	$\mathbf{w}_{2,opt} = \mathbf{u}_1\left(\left(eta^\star( ho)\mathbf{\Sigma}_1 + \mathbf{I} ight)^{-1}\mathbf{\Sigma}_2 ight)$
		$\{\alpha^{\star}(\rho), \beta^{\star}(\rho)\} \ge 0$ , chosen
		appropriately (See Conjecture 1)
$\rho \to \infty$	$\mathbf{w}_{i,opt} = \mathbf{u}_1\left(\mathbf{\Sigma}_i ight)$	$\mathbf{w}_{1,opt} = \mathbf{u}_1\left(\mathbf{\Sigma}_2^{-1}\mathbf{\Sigma}_1 ight)$
	$\mathbf{w}_{j,opt}=\mathbf{u}_{2}\left(\mathbf{\Sigma}_{i} ight),$	$\mathbf{w}_{2,opt} = \mathbf{u}_1\left(\mathbf{\Sigma}_1^{-1}\mathbf{\Sigma}_2 ight)$
	$j \neq i$ (See Prop. 4)	(See Theorems 2 and 3)

#### APPENDIX

## A. Density Function of Weighted-Norm of Isotropically Distributed Unit-Norm Vectors

Towards computing  $E[I_{i,1}]$ , we generalize the technique expounded in [51] where the surface area (that is required) to be computed is treated as a differential element of a corresponding solid volume (at a specific radius value), and the volume of the necessary solid object is calculated using tools from higher-dimensional integration (geometry). In this direction, we have

$$p_i(x) dx \triangleq \mathsf{P}\left(\widehat{\mathbf{h}}_{\mathsf{iid},i}^H \mathbf{\Lambda}_i \widehat{\mathbf{h}}_{\mathsf{iid},i} \in [x, x + dx]\right)$$
 (100)

$$p_i(x) = \frac{\partial}{\partial x} \mathsf{P}\left(\widehat{\mathbf{h}}_{\mathsf{iid},i}^H \mathbf{\Lambda}_i \widehat{\mathbf{h}}_{\mathsf{iid},i} \le x\right)$$
(101)

with

$$\mathsf{P}\left(\widehat{\mathbf{h}}_{\mathsf{iid},i}^{H} \mathbf{\Lambda}_{i} \widehat{\mathbf{h}}_{\mathsf{iid},i} \leq x\right) = 1 - \frac{\mathsf{Area}\left(x,\,1\right)}{\mathsf{Area}\left(1\right)} \tag{102}$$

where

Area 
$$(x, y) \triangleq$$
 Area  $\left(\widehat{\mathbf{h}}_{\mathsf{iid}, i}^{H} \mathbf{\Lambda}_{i} \widehat{\mathbf{h}}_{\mathsf{iid}, i} \ge x, \|\widehat{\mathbf{h}}_{\mathsf{iid}, i}\|^{2} = y\right)$  and (103)

Area 
$$(y) \triangleq$$
 Area  $\left( \|\widehat{\mathbf{h}}_{\mathsf{iid},i}\|^2 = y \right)$  (104)

denote the area of a (unit radius) spherical cap carved out by the ellipsoid

$$\left\{\widehat{\mathbf{h}}_{\mathsf{iid},i}:\widehat{\mathbf{h}}_{\mathsf{iid},i}^{H}\,\mathbf{\Lambda}_{i}\,\widehat{\mathbf{h}}_{\mathsf{iid},i}=x\right\}$$
(105)

and the area of a (unit radius) complex sphere, respectively. The volume of the objects desired in the computation of  $p_i(x)$  are

$$\mathsf{Vol}(x, r^2) \triangleq \mathsf{Vol}\left(\widehat{\mathbf{h}}_{\mathsf{iid}, i}^H \mathbf{\Lambda}_i \widehat{\mathbf{h}}_{\mathsf{iid}, i} \ge x, \|\widehat{\mathbf{h}}_{\mathsf{iid}, i}\|^2 \le r^2\right)$$
(106)

$$= \int_{y=0}^{r^{2}} \operatorname{Area}\left(x, y\right) dy, \tag{107}$$

$$\operatorname{Vol}(r^2) \triangleq \operatorname{Vol}\left(\|\widehat{\mathbf{h}}_{\operatorname{iid},i}\|^2 \le r^2\right) = \int_{x=0}^{r^2} \operatorname{Area}(x) dx.$$
(108)

Thus, we have

Area 
$$(x, 1) = \frac{\partial}{\partial r^2} \operatorname{Vol}(x, r^2) \Big|_{r=1},$$
 (109)

Area 
$$(x) = \frac{\partial}{\partial r^2} \operatorname{Vol}(r^2) \Big|_{r=1}$$
 and hence, (110)

$$p_i(x) = -\frac{\frac{\partial^2}{\partial x r^2} \operatorname{Vol}(x, r^2) \Big|_{r=1}}{\frac{\partial}{\partial r^2} \operatorname{Vol}(r^2) \Big|_{r=1}}.$$
(111)

It is important to realize that computing  $Vol(x, r^2)$  is *non-trivial* even in the simplest case of M = 2. This is because every additional dimension to the complex ellipsoid corresponds to addition of two real dimensions, thus rendering a geometric visualization impossible. For example, with M = 2, we have the intersection of two four-dimensional real objects. Nevertheless, the following lemma captures the complete structure of  $p_i(x)$  when M = 2.

*Lemma 1:* If M = 2, the random variable  $\widehat{\mathbf{h}}_{\text{iid}, i}^{H} \mathbf{\Lambda}_{i} \widehat{\mathbf{h}}_{\text{iid}, i}$  is uniformly distributed in the interval  $[\mathbf{\Lambda}_{i, 2}, \mathbf{\Lambda}_{i, 1}]$ .

Proof: First, note that it follows from [51, Lemma 2] that

=

$$Vol(r^{2}) = \frac{\pi^{M} r^{2M}}{M!}.$$
(112)

For computing Vol  $(x, r^2)$ , we follow the same variable transformation as in [51]. We set  $\widehat{\mathbf{h}}_{\text{iid},i}(k) = r_k \exp(j\phi_k)$  for k = 1, 2. The ellipsoid is contained completely in the sphere of radius r if r is such that  $r \ge \sqrt{\frac{x}{\Lambda_{i,2}}}$ , whereas the sphere is contained completely in the ellipsoid if  $r \le \sqrt{\frac{x}{\Lambda_{i,1}}}$ . In the intermediate regime for r, a non-trivial intersection between the two objects is observed and one can compute the volume by performing a two-dimensional integration as follows:

$$\operatorname{Vol}(x, r^2) = \iint_{\mathcal{A}} r_1 r_2 \phi_1 \phi_2 dr_1 dr_2 d\phi_1 d\phi_2 \tag{113}$$

$$= (2\pi)^2 \cdot \iint_{\mathcal{B}} r_1 dr_1 r_2 dr_2 \tag{114}$$

$$= (2\pi)^2 \cdot \int_0^{r^*} r_2 dr_2 \int_L^U r_1 dr_1$$
 (115)

where

$$\mathcal{A} = \left\{ r_1, r_2 : r_1^2 \Lambda_{i,1} + r_2^2 \Lambda_{i,2} \ge x, \ r_1^2 + r_2^2 \le r^2 \right\}$$
  
and  $\left\{ \phi_1, \phi_2 : [0, 2\pi) \right\},$  (116)

$$\mathcal{B} = \left\{ r_1, r_2 : r_1^2 \Lambda_{i,1} + r_2^2 \Lambda_{i,2} \ge x, \ r_1^2 + r_2^2 \le r^2 \right\},$$
(117)

$$L = \sqrt{\frac{x - r_2^2 \Lambda_{i,2}}{\Lambda_{i,1}}}, \quad U = \sqrt{r^2 - r_2^2}, \quad r^* = \frac{r^2 \Lambda_{i,1} - x}{\Lambda_{i,1} - \Lambda_{i,2}}.$$
 (118)

Trivial computation establishes the following:

$$\operatorname{Vol}\left(x,r^{2}\right) = \begin{cases} 0, & r \leq \sqrt{\frac{x}{\Lambda_{i,1}}} \\ \frac{\pi^{2}}{2} \cdot \frac{\left(r^{2}\Lambda_{i,1}-x\right)^{2}}{\Lambda_{i,1}\cdot\left(\Lambda_{i,1}-\Lambda_{i,2}\right)}, & \sqrt{\frac{x}{\Lambda_{i,1}}} \leq r \leq \sqrt{\frac{x}{\Lambda_{i,2}}} \\ \frac{\pi^{2}}{2} \cdot \left(r^{4}-\frac{x^{2}}{\Lambda_{i,1}\Lambda_{i,2}}\right), & r \geq \sqrt{\frac{x}{\Lambda_{i,2}}}. \end{cases}$$
(119)

Another trivial computation using (111) results in

$$p_i(x) = \frac{1}{\Lambda_{i,1} - \Lambda_{i,2}}.$$
(120)

That is,  $\widehat{\mathbf{h}}_{\text{iid},i}^{H} \mathbf{\Lambda}_{i} \widehat{\mathbf{h}}_{\text{iid},i}$  is uniformly distributed in its range. The structure of  $p_{i}(x)$  gets more complicated as M increases. We now provide its structure in the M = 3 and M = 4 cases without proof.

Lemma 2: With M = 3, the density function  $p_i(x)$  is of the form:

$$p_{i}(x) = \begin{cases} 0, & x \leq \Lambda_{i,3} \\ \frac{2(x - \Lambda_{i,3})}{(\Lambda_{i,1} - \Lambda_{i,3})(\Lambda_{i,2} - \Lambda_{i,3})}, & \Lambda_{i,3} \leq x \leq \Lambda_{i,2} \\ \frac{2(\Lambda_{i,1} - x)}{(\Lambda_{i,1} - \Lambda_{i,2})(\Lambda_{i,1} - \Lambda_{i,3})}, & \Lambda_{i,2} \leq x \leq \Lambda_{i,1} \\ 0, & x \geq \Lambda_{i,1}. \end{cases}$$
(121)

With M = 4, the density function  $p_i(x)$  takes the form:

$$p_{i}(x) = \begin{cases} 0, & x \leq \Lambda_{i,4} \\ \frac{3(x - \Lambda_{i,4})^{2}}{(\Lambda_{i,1} - \Lambda_{i,4})(\Lambda_{i,2} - \Lambda_{i,4})(\Lambda_{i,3} - \Lambda_{i,4})}, & \Lambda_{i,4} \leq x \leq \Lambda_{i,3} \\ \frac{3}{(\Lambda_{i,1} - \Lambda_{i,3})(\Lambda_{i,2} - \Lambda_{i,4})} \cdot \mathcal{L}_{0}, & \Lambda_{i,3} \leq x \leq \Lambda_{i,2} \\ \frac{3(\Lambda_{i,1} - x)^{2}}{(\Lambda_{i,1} - \Lambda_{i,2})(\Lambda_{i,1} - \Lambda_{i,3})(\Lambda_{i,1} - \Lambda_{i,4})}, & \Lambda_{i,2} \leq x \leq \Lambda_{i,1} \\ 0, & x \geq \Lambda_{i,1} \end{cases}$$
(122)

where

$$\mathcal{L}_{0} = \frac{(x - \Lambda_{i,3}) (\Lambda_{i,2} - x)}{\Lambda_{i,2} - \Lambda_{i,3}} + \frac{(x - \Lambda_{i,4}) (\Lambda_{i,1} - x)}{\Lambda_{i,1} - \Lambda_{i,4}}.$$
 (123)



Fig. 5. CDF of weighted-norm of isotropically distributed beamforming vectors.

Fig. 5 illustrates the trends of the cumulative distribution function (CDF) by plotting the fit between the theoretical expressions in Lemmas 1 and 2, and the CDF estimated by Monte Carlo methods. The cases considered are: a)  $\Lambda_i = \text{diag}([2 \ 1])$  for M = 2, b)  $\Lambda_i = \text{diag}([3 \ 2 \ 1])$  for M = 3, and c)  $\Lambda_i = \text{diag}([4 \ 3 \ 2 \ 1])$  for M = 4. The figure shows the excellent match between theory and Monte Carlo estimates.

#### B. Rewriting the Rate Expression in the High-SNR Extreme

A straightforward exercise shows that (36) can be rewritten as in (124) below:

$$E[R_i] \stackrel{\rho \to \infty}{\to} \frac{1}{2}g\left(d_{\Sigma_i}(\mathbf{w}_1, \mathbf{w}_2)\right) + \log\left(1 + \frac{A_i}{B_i}\right) - \log(2)$$
(124)

where  $g(\bullet)$  is a function defined as

$$g(z) \triangleq f(z) + 2\log(z), \tag{125}$$

$$f(z) = \frac{1}{\sqrt{1-z^2}} \log\left(\frac{1+\sqrt{1-z^2}}{1-\sqrt{1-z^2}}\right), \quad 0 < z < 1.$$
(126)

In (124),  $d_{\Sigma_i}(\mathbf{w}_1, \mathbf{w}_2)$  is defined as

$$d_{\Sigma_i} \left( \mathbf{w}_1, \mathbf{w}_2 \right) \triangleq \sqrt{\frac{4 \left( A_i B_i - C_i^2 \right)}{\left( A_i + B_i \right)^2}}$$
(127)

with  $A_i, B_i$  and  $C_i$  as in Theorem 1. As illustrated in Fig. 6,  $f(\bullet)$  is monotonically decreasing as a function of its argument and  $g(\bullet)$  is increasing with

$$2\log(2) = \lim_{z \to 0} g(z) \le g(z) \le \lim_{z \to 1} g(z) = 2$$
(128)

$$\infty = \lim_{z \to 0} f(z) \ge f(z) \ge \lim_{z \to 1} f(z) = 2.$$
(129)

Formal proofs of these facts are provided in Appendix C next. Some properties of  $d_{\Sigma_i}(\mathbf{w}_1, \mathbf{w}_2)$  are now established.

• The quantity  $d_{\Sigma_i}(\mathbf{w}_1, \mathbf{w}_2)$  is a generalized "distance" semi-metric<sup>16</sup> between  $\mathbf{w}_1$  and  $\mathbf{w}_2$  satisfying

$$0 \le d_{\Sigma_i} \left( \mathbf{w}_1, \mathbf{w}_2 \right) \le 1. \tag{130}$$

To establish the lower bound in (130), note that an application of the Cauchy-Schwarz inequality implies that  $C_i^2 \leq A_i B_i$ . Equality in the lower bound in (130) is achieved if and only if  $\Sigma_i^{1/2} \mathbf{w}_1 = \zeta \Sigma_i^{1/2} \mathbf{w}_2$  for some  $\zeta \in \mathbb{C}$ . Since  $\Sigma_i$  is positive-definite, this is possible only when  $\mathbf{w}_1 = \zeta \mathbf{w}_2$ . Since both  $\mathbf{w}_1$  and  $\mathbf{w}_2$  are unit-norm, this is possible only with  $|\zeta| = 1$ . In other words, equality in the lower bound only occurs for  $\mathbf{w}_1 = \mathbf{w}_2$  on  $\mathcal{G}(2, 1)$ .

<sup>&</sup>lt;sup>16</sup>A semi-metric satisfies all the properties necessary for a distance metric except the triangle inequality.



Fig. 6. The behavior of f(x) and g(x).

The fact that  $d_{\Sigma_i}(\mathbf{w}_1, \mathbf{w}_2) \leq 1$  is obvious. Symmetry of the distance metric in  $\mathbf{w}_1$  and  $\mathbf{w}_2$  is obvious.

• The triangle inequality does not hold in general. One counter-example is as follows:

$$\Sigma_i = \operatorname{diag}([20, 1]), \quad \mathbf{w}_1 = \left[\frac{1}{\sqrt{3}}, \sqrt{\frac{2}{3}}\right],$$
(131)

$$\mathbf{w}_2 = \left[\frac{1}{\sqrt{2}}, \frac{-1}{\sqrt{2}}\right], \quad \mathbf{w}_3 = \left[-\sqrt{\frac{3.3}{7}}, \sqrt{\frac{3.7}{7}}\right].$$
 (132)

This choice results in

$$d_{\Sigma_i}\left(\mathbf{w}_1, \mathbf{w}_3\right) \approx 0.2536, \quad d_{\Sigma_i}\left(\mathbf{w}_1, \mathbf{w}_2\right) + d_{\Sigma_i}\left(\mathbf{w}_2, \mathbf{w}_3\right) \approx 0.2534.$$
(133)

Many such counter-examples can be listed out via a routine numerical search. Numerical studies also suggest that the triangle inequality holds for almost all choices of  $\{\mathbf{w}_i\}$  provided that  $\chi_i = \frac{\lambda_1(\Sigma_i)}{\lambda_2(\Sigma_i)}$  is not too large (unlike the example in (131)-(132)).

The upper bound in (130) is achieved only if A<sub>i</sub> = B<sub>i</sub> and C<sub>i</sub> = 0. By decomposing w<sub>1</sub> and w<sub>2</sub> along the orthonormal set of basis vectors {u<sub>1</sub>(Σ<sub>i</sub>), u<sub>2</sub>(Σ<sub>i</sub>)}, it can be checked that A<sub>i</sub> = B<sub>i</sub> and C<sub>i</sub> = 0 is possible only if

$$\mathbf{w}_1 = e^{j\nu_1} \cdot \left[ \mathbf{u}_1(\mathbf{\Sigma}_i) \cdot \sqrt{\frac{1}{\chi_i + 1}} + e^{j\nu_2} \cdot \mathbf{u}_2(\mathbf{\Sigma}_i) \cdot \sqrt{\frac{\chi_i}{\chi_i + 1}} \right]$$
(134)

$$\mathbf{w}_2 = e^{j(\nu_1 + \nu_3)} \cdot \left[ \mathbf{u}_1(\boldsymbol{\Sigma}_i) \cdot \sqrt{\frac{1}{\chi_i + 1}} - e^{j\nu_2} \cdot \mathbf{u}_2(\boldsymbol{\Sigma}_i) \cdot \sqrt{\frac{\chi_i}{\chi_i + 1}} \right]$$
(135)

for some choice of  $\nu_j \in [0, 2\pi), j = 1, 2, 3$ .

• The semi-metric reduces to the standard *chordal* distance metric [42] on  $\mathcal{G}(2,1)$ 

$$d_{\mathbf{\Sigma}_i}\left(\mathbf{w}_1, \mathbf{w}_2\right) = \sqrt{1 - |\mathbf{w}_1^H \mathbf{w}_2|^2}$$
(136)

if  $\Sigma_i = \lambda \mathbf{I}$  for some  $\lambda > 0$ .

C. Monotonicity of  $f(\bullet)$  and  $g(\bullet)$ 

We first claim that

$$2 \le \frac{1}{\sqrt{1-z^2}} \cdot \log\left(\frac{1+\sqrt{1-z^2}}{1-\sqrt{1-z^2}}\right) \le \frac{2}{z^2}, \quad 0 < z < 1,$$
(137)

which is equivalent to:

$$\exp\left(2\cdot\sqrt{1-z^2}\right) \le \frac{1+\sqrt{1-z^2}}{1-\sqrt{1-z^2}} \le \exp\left(\frac{2\sqrt{1-z^2}}{z^2}\right), \quad 0 < z < 1.$$
(138)

For this, we start with the exponential series expansion of  $\exp\left(2\cdot\sqrt{1-z^2}\right)$  that results in:

$$\exp\left(2\cdot\sqrt{1-z^2}\right) = 1 + \sum_{k=1}^{\infty} \frac{2^k \cdot (1-z^2)^{\frac{k}{2}}}{\Gamma(k+1)}$$
(139)

$$\leq 1 + 2\sum_{k=1}^{\infty} \left(1 - z^2\right)^{\frac{k}{2}} = 1 + \frac{2\sqrt{1 - z^2}}{1 - \sqrt{1 - z^2}}$$
(140)

where  $\Gamma(\cdot)$  is the Gamma function, the second inequality follows from the fact that  $\frac{2^{k-1}}{\Gamma(k+1)} \leq 1$  for all  $k \geq 1$ , and the last equality from the sum of an infinite geometric series. For the other side of (137), note that

$$\exp\left(\frac{2\sqrt{1-z^2}}{z^2}\right) \ge 1 + \frac{2\sqrt{1-z^2}}{z^2} + \frac{2(1-z^2)}{z^4}$$
(141)

$$\geq 1 + \frac{2\sqrt{1-z^2}}{z^2} \left(1 + \sqrt{1-z^2}\right) = 1 + \frac{2\sqrt{1-z^2}}{1 - \sqrt{1-z^2}}$$
(142)

where the first inequality follows by truncating the terms of the asymptotic expansion and the second follows by using the fact that  $z^2 < 1$ . The proof is complete by noting that

$$\frac{\partial f(z)}{\partial z} = \frac{-z}{1-z^2} \cdot \left[\frac{2}{z^2} - \frac{1}{\sqrt{1-z^2}} \log\left(\frac{1+\sqrt{1-z^2}}{1-\sqrt{1-z^2}}\right)\right] < 0$$
(143)

$$\frac{\partial g(z)}{\partial z} = \frac{z}{1-z^2} \cdot \left[ \frac{1}{\sqrt{1-z^2}} \log\left(\frac{1+\sqrt{1-z^2}}{1-\sqrt{1-z^2}}\right) - 2 \right] > 0.$$
(144)

## D. Completing the Proof of Theorem 2

With the description of  $w_1$  and  $w_2$  as in (48)-(49), elementary algebra shows that

$$A_{1} = \mathbf{w}_{1}^{H} \boldsymbol{\Sigma}_{1} \mathbf{w}_{1} = \frac{|\alpha|^{2} \eta_{1} + |\beta|^{2} \eta_{2}}{X^{2}}, \ X = \|\alpha \boldsymbol{\Sigma}_{2}^{-\frac{1}{2}} \mathbf{v}_{1} + \beta \boldsymbol{\Sigma}_{2}^{-\frac{1}{2}} \mathbf{v}_{2}\|$$
(145)

$$B_{1} = \mathbf{w}_{2}^{H} \boldsymbol{\Sigma}_{1} \mathbf{w}_{2} = \frac{|\gamma|^{2} \eta_{1} + |\delta|^{2} \eta_{2}}{Y^{2}}, \quad Y = \|\gamma \boldsymbol{\Sigma}_{2}^{-\frac{1}{2}} \mathbf{v}_{1} + \delta \boldsymbol{\Sigma}_{2}^{-\frac{1}{2}} \mathbf{v}_{2}\|$$
(146)

$$C_1 = |\mathbf{w}_1^H \mathbf{\Sigma}_1 \mathbf{w}_2| = \frac{\left|\alpha^* \gamma \eta_1 + \beta^* \delta \eta_2\right|}{XY}$$
(147)

$$A_2 = \mathbf{w}_2^H \boldsymbol{\Sigma}_2 \mathbf{w}_2 = \frac{1}{Y^2}$$
(148)

$$B_2 = \mathbf{w}_1^H \boldsymbol{\Sigma}_2 \mathbf{w}_1 = \frac{1}{X^2}$$
(149)

$$C_2 = |\mathbf{w}_1^H \boldsymbol{\Sigma}_2 \mathbf{w}_2| = \frac{\left|\alpha^* \gamma + \beta^* \delta\right|}{XY}.$$
(150)

As in the statement of the theorem, let  $\tau_i$ , i = 1, 2, 3 denote

$$\tau_1 = \mathbf{v}_1^H \boldsymbol{\Sigma}_2^{-1} \mathbf{v}_1, \quad \tau_2 = \mathbf{v}_2^H \boldsymbol{\Sigma}_2^{-1} \mathbf{v}_2, \quad \tau_3 = \mathbf{v}_1^H \boldsymbol{\Sigma}_2^{-1} \mathbf{v}_2.$$
(151)

We can rewrite  $X^2$  and  $Y^2$  in terms of  $\{\tau_i\}$  as

$$X^{2} = |\alpha|^{2}\tau_{1} + |\beta|^{2}\tau_{2} + 2|\alpha||\beta||\tau_{3}|\cos(\theta_{1})$$
(152)

$$Y^{2} = |\gamma|^{2} \tau_{1} + |\delta|^{2} \tau_{2} + 2|\gamma| |\delta| |\tau_{3}| \cos(\theta_{2})$$
(153)

where  $\theta_1 = \arg(\tau_3) + \theta_\beta - \theta_\alpha$  and  $\theta_2 = \arg(\tau_3) + \theta_\delta - \theta_\gamma$ . Now note that if  $\{\mathbf{v}_1, \mathbf{v}_2\}$  is a pair of eigenvectors for  $\Sigma$ , then so is the pair  $\{e^{j\nu_1}\mathbf{v}_1, e^{j\nu_2}\mathbf{v}_2\}$  for any choice of  $\nu_1$  and  $\nu_2$ . In other words, the choice of  $\{\mathbf{v}_1, \mathbf{v}_2\}$  is unique *only* on  $\mathcal{G}(2, 1)$ , and not on  $\mathsf{St}(2, 1)$ . Hence,  $\arg(\tau_3)$ can be chosen arbitrarily and *independently* in determining the values of  $X^2$  and  $Y^2$ . With the specific choice that  $\arg(\mathbf{v}_1^H \Sigma_2^{-1} \mathbf{v}_2) = \frac{\pi}{2} + \theta_\alpha - \theta_\beta$  in (152) and  $\arg(\mathbf{v}_1^H \Sigma_2^{-1} \mathbf{v}_2) = \frac{\pi}{2} + \theta_\gamma - \theta_\delta$ in (153), we have

$$X^{2} = |\alpha|^{2}\tau_{1} + |\beta|^{2}\tau_{2}$$
(154)

$$Y^{2} = |\gamma|^{2} \tau_{1} + |\delta|^{2} \tau_{2}.$$
(155)

Thus, the high-SNR expression for the ergodic sum-rate can be simplified as

$$2E[R_{1}] + 2E[R_{2}] + 4\log(2) = g\left(\frac{2\sqrt{\eta_{1}\eta_{2}}XY \cdot |\beta\gamma - \alpha\delta|}{\left(|\alpha|^{2}\eta_{1} + |\beta|^{2}\eta_{2}\right) \cdot Y^{2} + \left(|\gamma|^{2}\eta_{1} + |\delta|^{2}\eta_{2}\right) \cdot X^{2}}\right) + 2\log\left(1 + \frac{X^{2}}{Y^{2}}\right) + g\left(\frac{2XY \cdot |\beta\gamma - \alpha\delta|}{X^{2} + Y^{2}}\right) + 2\log\left(1 + \frac{Y^{2}}{X^{2}} \cdot \frac{|\alpha|^{2}\eta_{1} + |\beta|^{2}\eta_{2}}{|\gamma|^{2}\eta_{1} + |\delta|^{2}\eta_{2}}\right).$$
(156)

Using (126) to rewrite the above equation in terms of  $f(\bullet)$ , we have after simplification:

$$2E[R_{1}] + 2E[R_{2}] - \log(\eta_{1}\eta_{2}) = f\left(\frac{2\sqrt{\eta_{1}\eta_{2}}XY \cdot |\beta\gamma - \alpha\delta|}{\left(|\alpha|^{2}\eta_{1} + |\beta|^{2}\eta_{2}\right)Y^{2} + \left(|\gamma|^{2}\eta_{1} + |\delta|^{2}\eta_{2}\right)X^{2}}\right) + f\left(\frac{2XY \cdot |\beta\gamma - \alpha\delta|}{X^{2} + Y^{2}}\right) + 2\log\left(\frac{|\beta\gamma - \alpha\delta|^{2}}{|\gamma|^{2}\eta_{1} + |\delta|^{2}\eta_{2}}\right).$$
(157)

We now claim that the following two inequalities hold:

$$\frac{XY}{X^2 + Y^2} \ge \frac{\sqrt{\tau_1 \tau_2}}{\tau_1 + \tau_2}$$
(158)

$$\frac{XY \cdot \left(|\gamma|^2 \eta_1 + |\delta|^2 \eta_2\right)}{\left(|\alpha|^2 \eta_1 + |\beta|^2 \eta_2\right)Y^2 + \left(|\gamma|^2 \eta_1 + |\delta|^2 \eta_2\right)X^2} \geq \frac{\sqrt{\tau_1 \tau_2} \cdot \eta_2}{\tau_1 \eta_2 + \tau_2 \eta_1}.$$
(159)

The proof of (158) and (159) will be tackled later.

Using (158) and (159) in conjunction with the decreasing nature of  $f(\bullet)$ , we have

$$2E[R_1] + 2E[R_2] - \log(\eta_1\eta_2) \leq f\left(\frac{2\sqrt{\eta_1\eta_2\tau_1\tau_2}\cdot\eta_2}{\tau_1\eta_2+\tau_2\eta_1}\cdot\frac{|\beta\gamma-\alpha\delta|}{(|\gamma|^2\eta_1+|\delta|^2\eta_2)}\right) + f\left(\frac{2\sqrt{\tau_1\tau_2}}{\tau_1+\tau_2}\cdot|\beta\gamma-\alpha\delta|\right) + 2\log\left(\frac{|\beta\gamma-\alpha\delta|^2}{|\gamma|^2\eta_1+|\delta|^2\eta_2}\right)$$
(160)  
$$= g\left(\frac{2\sqrt{\eta_1\eta_2\tau_1\tau_2}\cdot\eta_2}{\tau_1\eta_2+\tau_2\eta_1}\cdot\frac{|\beta\gamma-\alpha\delta|}{(|\gamma|^2\eta_1+|\delta|^2\eta_2)}\right) + g\left(\frac{2\sqrt{\tau_1\tau_2}}{\tau_1+\tau_2}\cdot|\beta\gamma-\alpha\delta|\right) + 2\log\left(\frac{(\tau_1\eta_2+\tau_2\eta_1)(\tau_1+\tau_2)}{4\tau_1\tau_2\eta_2\sqrt{\eta_1\eta_2}}\right).$$
(161)

Note that  $\{\theta_{\bullet}\}$  enter the above optimization only via the term  $|\beta\gamma-\alpha\delta|$  and

$$|\beta\gamma - \alpha\delta| \le |\alpha|\sqrt{1 - |\gamma|^2} + |\gamma|\sqrt{1 - |\alpha|^2}$$
(162)

with equality achieved if and only if  $\theta_{\alpha} + \theta_{\delta} - \theta_{\beta} - \theta_{\gamma} = \pi$  (modulo  $2\pi$ ). Parameterizing  $|\alpha|$  and  $|\gamma|$  as  $|\alpha| = \sin(\theta)$  and  $|\gamma| = \sin(\phi)$  for some  $\{\theta, \phi\} \in [0, \pi/2]$ , we have

$$|\beta\gamma - \alpha\delta| \le \sin(\theta)\cos(\phi) + \cos(\theta)\sin(\phi) = \sin(\theta + \phi) \le 1$$
(163)

since  $0 \le \theta + \phi \le \pi$ . Further,  $\eta_1 \ge \eta_2$  implies that  $|\gamma|^2 \eta_1 + |\delta|^2 \eta_2 \ge \eta_2$  and hence, we have

$$\frac{|\beta\gamma - \alpha\delta|}{|\gamma|^2\eta_1 + |\delta|^2\eta_2} \le \frac{1}{\eta_2}.$$
(164)

Using the fact that  $g(\bullet)$  is an increasing function, we get an upper bound for the sum-rate as

$$2E[R_1] + 2E[R_2] \le f\left(\frac{2\sqrt{\eta_1\eta_2\tau_1\tau_2}}{\eta_1\tau_2 + \eta_2\tau_1}\right) + f\left(\frac{2\sqrt{\tau_1\tau_2}}{\tau_1 + \tau_2}\right) + \log\left(\frac{\eta_1}{\eta_2}\right).$$
(165)

This upper bound is achievable with the choice of  $|\alpha| = 1$  and  $|\gamma| = 0$  in (48) and (49), which is equivalent to (44). Substituting this choice in the sum-rate expression yields

$$E[R_1] + E[R_2] \xrightarrow{\rho \to \infty} \frac{1}{2} f\left(\frac{2\sqrt{\eta_1 \eta_2 \tau_1 \tau_2}}{\eta_1 \tau_2 + \eta_2 \tau_1}\right) + \frac{1}{2} f\left(\frac{2\sqrt{\tau_1 \tau_2}}{\tau_1 + \tau_2}\right) + \frac{1}{2} \log\left(\frac{\eta_1}{\eta_2}\right)$$
(166)

$$= \frac{1}{2} \cdot \frac{\eta_{1}\tau_{2} + \eta_{2}\tau_{1}}{|\eta_{1}\tau_{2} - \eta_{2}\tau_{1}|} \cdot \log\left(\frac{\eta_{1}\tau_{2} + \eta_{2}\tau_{1} + |\eta_{1}\tau_{2} - \eta_{2}\tau_{1}|}{\eta_{1}\tau_{2} + \eta_{2}\tau_{1} - |\eta_{1}\tau_{2} - \eta_{2}\tau_{1}|}\right) \\ + \frac{1}{2} \cdot \frac{\tau_{1} + \tau_{2}}{|\tau_{1} - \tau_{2}|} \cdot \log\left(\frac{\tau_{1} + \tau_{2} + |\tau_{1} - \tau_{2}|}{\tau_{1} + \tau_{2} - |\tau_{1} - \tau_{2}|}\right) + \frac{1}{2}\log\left(\frac{\eta_{1}}{\eta_{2}}\right).$$
(167)

To simplify (167), we define  $\kappa_1$  and  $\kappa_2$  as

$$\kappa_1 = \frac{\eta_1 \tau_2}{\eta_2 \tau_1} \quad \text{and} \quad \kappa_2 = \frac{\tau_2}{\tau_1}.$$
(168)

The fact that  $\eta_1 \ge \eta_2$  implies that  $\kappa_1 \ge \kappa_2$ . Thus, there are three possibilities: i)  $\kappa_1 \ge \kappa_2 \ge 1$ , ii)  $\kappa_1 \ge 1 \ge \kappa_2$ , and iii)  $1 \ge \kappa_1 \ge \kappa_2$ . It is straightforward but tedious to check that in all of the three cases, (167) reduces to (45) as in the statement of the theorem. The proof will be complete if the inequalities (158) and (159) can be established.

**Proof of (158):** For the first inequality, note that

$$\frac{X^2 + Y^2}{XY} = \frac{X}{Y} + \frac{Y}{X} = t + \frac{1}{t} \triangleq q(t)$$
(169)

can be written as a symmetric function q(t) in t where  $t = \frac{X}{Y}$ . Further, noting that q(t) is decreasing in t for  $t \le 1$  and is increasing in t for  $t \ge 1$ , the maximum of  $\frac{X^2+Y^2}{XY}$  is achieved either when  $\frac{X}{Y}$  achieves its largest or smallest value. The inequality in (158) follows since

$$\sqrt{\frac{\min(\tau_1, \tau_2)}{\max(\tau_1, \tau_2)}} \le \frac{X}{Y} \le \sqrt{\frac{\max(\tau_1, \tau_2)}{\min(\tau_1, \tau_2)}}.$$
(170)

Proof of (159): The proof of (159) is more involved. For this, note that

$$\mathcal{L}_{1} \triangleq \frac{\left(|\alpha|^{2}\eta_{1} + |\beta|^{2}\eta_{2}\right)Y^{2} + \left(|\gamma|^{2}\eta_{1} + |\delta|^{2}\eta_{2}\right)X^{2}}{XY \cdot \left(|\gamma|^{2}\eta_{1} + |\delta|^{2}\eta_{2}\right)}$$
(171)

$$= \sqrt{\frac{|\alpha|^2(\tau_1 - \tau_2) + \tau_2}{|\gamma|^2(\tau_1 - \tau_2) + \tau_2}} \cdot \left(1 + \frac{|\alpha|^2(\eta_1 - \eta_2) + \eta_2}{|\gamma|^2(\eta_1 - \eta_2) + \eta_2} \cdot \frac{|\gamma|^2(\tau_1 - \tau_2) + \tau_2}{|\alpha|^2(\tau_1 - \tau_2) + \tau_2}\right).$$
(172)

By taking derivative with respect to  $|\alpha|^2$ , note that the first term in (172) is increasing in  $|\alpha|^2$  for any fixed choice of  $|\gamma|$  if and only if  $\frac{\tau_1}{\tau_2} \ge 1$ . Similarly, for any fixed choice of  $|\gamma|$ , the second term in (172) is increasing in  $|\alpha|^2$  if and only if  $\frac{\eta_1}{\eta_2} \ge \frac{\tau_1}{\tau_2}$ . Thus, the condition

$$1 \le \frac{\tau_1}{\tau_2} \le \frac{\eta_1}{\eta_2} \tag{173}$$

is necessary and sufficient to ensure that for any choice of  $|\gamma|$ ,  $\mathcal{L}_1$  is maximized by the choice  $|\alpha| = 1$ . On analogous lines, taking the derivative with respect to  $|\gamma|^2$ , it can be seen that for any fixed choice of  $|\alpha|$ , both the terms in (172) are decreasing in  $|\gamma|^2$  if and only if the same

condition in (173) holds. In other words, under (173),  $\mathcal{L}_1$  is maximized by  $|\alpha| = 1$  and  $|\gamma| = 0$ . At this stage, two other possibilities need to be considered: i)  $\frac{\tau_1}{\tau_2} \leq 1 \leq \frac{\eta_1}{\eta_2}$ , and ii)  $1 \leq \frac{\eta_1}{\eta_2} \leq \frac{\tau_1}{\tau_2}$ . In either case, we will show that

$$\left(\frac{|\alpha|^2(\eta_1 - \eta_2) + \eta_2}{|\gamma|^2(\eta_1 - \eta_2) + \eta_2} - \frac{\eta_1}{\eta_2}\right) \cdot \frac{Y}{X} + \left(\frac{X}{Y} - \sqrt{\frac{\tau_1}{\tau_2}}\right) \cdot \left(1 - \frac{\eta_1}{\eta_2}\sqrt{\frac{\tau_2}{\tau_1}}\frac{Y}{X}\right) \le 0,$$
(174)

which is equivalent to (159), or the statement that  $\mathcal{L}_1$  is maximized by  $|\alpha| = 1$  and  $|\gamma| = 0$ . For this, note that in either case, we have

$$\frac{|\alpha|^2(\eta_1 - \eta_2) + \eta_2}{|\gamma|^2(\eta_1 - \eta_2) + \eta_2} \le \frac{\eta_1}{\eta_2}.$$
(175)

In the first case, we also have

$$\sqrt{\frac{\tau_1}{\tau_2}} \le \frac{X}{Y} \le \sqrt{\frac{\tau_2}{\tau_1}} \le \frac{\eta_1}{\eta_2} \sqrt{\frac{\tau_2}{\tau_1}},\tag{176}$$

where the last step in (176) follows from  $\frac{\eta_1}{\eta_2} \ge 1$ . Combining (175) and (176), we note that (174) is immediate when  $\frac{\tau_1}{\tau_2} \le 1 \le \frac{\eta_1}{\eta_2}$ . In the second case, however, (176) is replaced with

$$\sqrt{\frac{\tau_2}{\tau_1}} \le \frac{X}{Y} \le \sqrt{\frac{\tau_1}{\tau_2}}.$$
(177)

It can be seen that if  $|\alpha|$  and  $|\gamma|$  are such that

$$\frac{X}{Y} = D \cdot \frac{\eta_1}{\eta_2} \sqrt{\frac{\tau_2}{\tau_1}},\tag{178}$$

for some choice of D satisfying  $1 \le D \le \frac{\tau_1}{\tau_2} \cdot \frac{\eta_2}{\eta_1}$ , (174) holds immediately. Thus, we only need to show that (174) holds when  $|\alpha|$  and  $|\gamma|$  are such that

$$\frac{X}{Y} = D \cdot \frac{\eta_1}{\eta_2} \sqrt{\frac{\tau_2}{\tau_1}},\tag{179}$$

for some choice of D satisfying  $\frac{\eta_2}{\eta_1} \le D \le 1$ . After some elementary algebra, our task is to show that

$$\frac{|\alpha|^2(\eta_1 - \eta_2) + \eta_2}{|\gamma|^2(\eta_1 - \eta_2) + \eta_2} \leq \frac{\eta_1}{\eta_2} \cdot \left(1 - (1 - D) \cdot \left(1 - D \cdot \frac{\eta_1}{\eta_2} \cdot \frac{\tau_2}{\tau_1}\right)\right)$$
(180)

$$\frac{X}{Y} = D \cdot \frac{\eta_1}{\eta_2} \sqrt{\frac{\tau_2}{\tau_1}}$$
(181)

By bounding the denominator of (180) as  $\eta_2 \leq |\gamma|^2(\eta_1 - \eta_2) + \eta_2 \leq \eta_1$ , it can be seen that (180) holds if the following quadratic inequality in D is true:

$$D^{2} \cdot \frac{\eta_{1}^{2} \tau_{2}}{\eta_{2}^{2} \tau_{1}} \cdot \left(2 - \frac{\tau_{1} \eta_{2} - \tau_{2} \eta_{1}}{\eta_{1} (\tau_{1} - \tau_{2})}\right) - D \cdot \frac{\eta_{1}}{\eta_{2}} \cdot \left(1 + \frac{\eta_{1} \tau_{2}}{\eta_{2} \tau_{1}}\right) + \frac{\tau_{1} \eta_{2} - \tau_{2} \eta_{1}}{\eta_{2} (\tau_{1} - \tau_{2})} \leq 0.$$
(182)

For this, we note that the left-hand side represents a convex parabola in D with maximum achieved at either  $D = \frac{\eta_2}{\eta_1}$  or D = 1. Substituting  $D = \frac{\eta_2}{\eta_1}$  and D = 1 and simplifying, we see that

LHS of (182)
$$\Big|_{D=\frac{\eta_2}{\eta_1}} = -\frac{\tau_2}{\eta_1\eta_2\tau_1} \cdot \frac{\eta_1 - \eta_2}{\tau_1 - \tau_2} \cdot (\eta_1(\tau_1 - \tau_2) + \tau_1(\eta_1 - \eta_2)) \le 0$$
 (183)

LHS of (182)
$$\Big|_{D=1} = -\frac{(\eta_1 - \eta_2) \cdot (\eta_2 \tau_1 - \eta_1 \tau_2)}{\eta_2^2 (\tau_1 - \tau_2)} \le 0.$$
 (184)

Since the maximum of the parabola in the domain  $\frac{\eta_2}{\eta_1} \le D \le 1$  is below 0, (174) holds. Thus, we are done with the aspect of showing that  $|\alpha| = 1$ ,  $|\gamma| = 0$  is sum-rate optimal.

## E. Proof of Theorem 3

Following the logic of Appendix D, we decompose  $\mathbf{w}_1$  and  $\mathbf{w}_2$  along  $\{\mathbf{u}_1, \mathbf{u}_2\}$  since they form an orthonormal basis:

$$\mathbf{w}_1 = \alpha \mathbf{u}_1 + \beta \mathbf{u}_2 \tag{185}$$

$$\mathbf{w}_2 = \gamma \mathbf{u}_1 + \delta \mathbf{u}_2 \tag{186}$$

for some choice of  $\{\alpha, \beta, \gamma, \delta\}$  with  $\alpha = |\alpha|e^{j\theta_{\alpha}}$  (similarly, for other quantities) satisfying  $|\alpha|^2 + |\beta|^2 = |\gamma|^2 + |\delta|^2 = 1$ . We now study the ergodic sum-rate optimization over the six-dimensional parameter space. With the description of  $\mathbf{w}_1$  and  $\mathbf{w}_2$  as in (185)-(186), elementary algebra shows that

$$A_1 = \mathbf{w}_1^H \boldsymbol{\Sigma}_1 \mathbf{w}_1 = |\alpha|^2 \lambda_1 + |\beta|^2 \lambda_2$$
(187)

$$B_1 = \mathbf{w}_2^H \boldsymbol{\Sigma}_1 \mathbf{w}_2 = |\gamma|^2 \lambda_1 + |\delta|^2 \lambda_2$$
(188)

$$C_1 = |\mathbf{w}_1^H \mathbf{\Sigma}_1 \mathbf{w}_2| = |\alpha^* \gamma \lambda_1 + \beta^* \delta \lambda_2|$$
(189)

$$A_2 = \mathbf{w}_2^H \mathbf{\Sigma}_2 \mathbf{w}_2 = |\gamma|^2 \mu_1 + |\delta|^2 \mu_2$$
(190)

$$B_2 = \mathbf{w}_1^H \mathbf{\Sigma}_2 \mathbf{w}_1 = |\alpha|^2 \mu_1 + |\beta|^2 \mu_2$$
(191)

$$C_2 = |\mathbf{w}_1^H \boldsymbol{\Sigma}_2 \mathbf{w}_2| = |\alpha^* \gamma \mu_1 + \beta^* \delta \mu_2|$$
(192)

and hence,

$$d_{\Sigma_1}(\mathbf{w}_1, \mathbf{w}_2)^2 = \frac{4(A_1 B_1 - C_1^2)}{(A_1 + B_1)^2} = \frac{4\lambda_1 \lambda_2 \cdot |\beta\gamma - \alpha\delta|^2}{\left[(|\alpha|^2 + |\gamma|^2)\lambda_1 + (|\beta|^2 + |\delta|^2)\lambda_2\right]^2}$$
(193)

$$d_{\Sigma_2}(\mathbf{w}_1, \mathbf{w}_2)^2 = \frac{4(A_2B_2 - C_2^2)}{(A_2 + B_2)^2} = \frac{4\mu_1\mu_2 \cdot |\beta\gamma - \alpha\delta|^2}{\left[(|\alpha|^2 + |\gamma|^2)\mu_1 + (|\beta|^2 + |\delta|^2)\mu_2\right]^2}.$$
 (194)

The high-SNR expression of the ergodic sum-rate can be written as

$$2E[R_{1}] + 2E[R_{2}] + 4\log(2) = g\left(\frac{\sqrt{4\lambda_{1}\lambda_{2}} \cdot |\beta\gamma - \alpha\delta|}{(|\alpha|^{2} + |\gamma|^{2})\lambda_{1} + (|\beta|^{2} + |\delta|^{2})\lambda_{2}}\right) + g\left(\frac{\sqrt{4\mu_{1}\mu_{2}} \cdot |\beta\gamma - \alpha\delta|}{(|\alpha|^{2} + |\gamma|^{2})\mu_{1} + (|\beta|^{2} + |\delta|^{2})\mu_{2}}\right) + 2\log\left(1 + \frac{|\alpha|^{2}\lambda_{1} + |\beta|^{2}\lambda_{2}}{|\gamma|^{2}\lambda_{1} + |\delta|^{2}\lambda_{2}}\right) + 2\log\left(1 + \frac{|\gamma|^{2}\mu_{1} + |\delta|^{2}\mu_{2}}{|\alpha|^{2}\mu_{1} + |\beta|^{2}\mu_{2}}\right).$$
(195)

Note that  $\{\theta_{\bullet}\}$  enter the above optimization only via the term  $|\beta\gamma - \alpha\delta|$  and as in (162), we have

$$|\beta\gamma - \alpha\delta| \le |\beta||\gamma| + |\alpha||\delta| \le 1.$$
(196)

Given that  $\chi_1 = \frac{\lambda_1}{\lambda_2} \ge 1$ , three possibilities arise depending on the relationship between 1,  $\chi_1$  and  $\chi_2 = \frac{\mu_1}{\mu_2}$ : i)  $\chi_1 > 1 \ge \chi_2$ , ii)  $\chi_1 > \chi_2 > 1$ , and iii)  $\chi_2 \ge \chi_1 > 1$ .

*Case i*): In the first case where  $\chi_2 \leq 1$ , we use the fact that  $g(\bullet)$  is an increasing function to bound the sum-rate as

$$2E[R_1] + 2E[R_2] + 4\log(2)$$

$$\leq f\left(\frac{\sqrt{4\lambda_1\lambda_2} \cdot \sin(\theta + \phi)}{(\sin^2(\theta) + \sin^2(\phi))(\lambda_1 - \lambda_2) + 2\lambda_2}\right) + f\left(\frac{\sqrt{4\mu_1\mu_2} \cdot \sin(\theta + \phi)}{2\mu_2 - (\sin^2(\theta) + \sin^2(\phi))(\mu_2 - \mu_1)}\right)$$

$$+ 2\log\left(\frac{\sqrt{4\lambda_1\lambda_2} \cdot \sin(\theta + \phi)}{\sin^2(\phi)(\lambda_1 - \lambda_2) + \lambda_2}\right) + 2\log\left(\frac{\sqrt{4\mu_1\mu_2} \cdot \sin(\theta + \phi)}{\mu_2 - \sin^2(\theta)(\mu_2 - \mu_1)}\right).$$
(197)

Observing that

$$(\sin^2(\theta) + \sin^2(\phi))(\lambda_1 - \lambda_2) + 2\lambda_2 \leq \sin^2(\phi)(\lambda_1 - \lambda_2) + \lambda_1 + \lambda_2$$
(198)

$$2\mu_2 - (\sin^2(\theta) + \sin^2(\phi))(\mu_2 - \mu_1) \le 2\mu_2 - \sin^2(\theta)(\mu_2 - \mu_1)$$
(199)

and  $f(\bullet)$  is a decreasing function, we have

$$2E[R_{1}] + 2E[R_{2}] + 4\log(2) \\ \leq f\left(\frac{\sqrt{4\lambda_{1}\lambda_{2}} \cdot \sin(\theta + \phi)}{\sin^{2}(\phi)(\lambda_{1} - \lambda_{2}) + \lambda_{1} + \lambda_{2}}\right) + f\left(\frac{\sqrt{4\mu_{1}\mu_{2}} \cdot \sin(\theta + \phi)}{2\mu_{2} - \sin^{2}(\theta)(\mu_{2} - \mu_{1})}\right) \\ + 2\log\left(\frac{\sqrt{4\lambda_{1}\lambda_{2}} \cdot \sin(\theta + \phi)}{\sin^{2}(\phi)(\lambda_{1} - \lambda_{2}) + \lambda_{2}}\right) + 2\log\left(\frac{\sqrt{4\mu_{1}\mu_{2}} \cdot \sin(\theta + \phi)}{\mu_{2} - \sin^{2}(\theta)(\mu_{2} - \mu_{1})}\right) \\ = g\left(\frac{\sqrt{4\lambda_{1}\lambda_{2}} \cdot \sin(\theta + \phi)}{\sin^{2}(\phi)(\lambda_{1} - \lambda_{2}) + \lambda_{1} + \lambda_{2}}\right) + g\left(\frac{\sqrt{4\mu_{1}\mu_{2}} \cdot \sin(\theta + \phi)}{2\mu_{2} - \sin^{2}(\theta)(\mu_{2} - \mu_{1})}\right) \\ + 2\log\left(1 + \frac{\lambda_{1}}{\sin^{2}(\phi)(\lambda_{1} - \lambda_{2}) + \lambda_{2}}\right) + 2\log\left(1 + \frac{\mu_{2}}{\mu_{2} - \sin^{2}(\theta)(\mu_{2} - \mu_{1})}\right).$$
(201)

It is straightforward to note that the right-hand side of (201) is decreasing in  $\sin^2(\phi)$  and increasing in  $\sin^2(\theta)$ . If in addition, the condition that  $\theta + \phi = \pi/2$  is satisfied, then an upper

bound on  $E[R_1] + E[R_2]$  can be maximized. This results in the choice  $\theta = \pi/2$  and  $\phi = 0$  (that is,  $|\alpha| = 1$  and  $|\gamma| = 0$ ) and with this choice, we have

$$2E[R_1] + 2E[R_2] \le f\left(\frac{\sqrt{4\lambda_1\lambda_2}}{\lambda_1 + \lambda_2}\right) + f\left(\frac{\sqrt{4\mu_1\mu_2}}{\mu_1 + \mu_2}\right) + \log\left(\frac{\lambda_1}{\lambda_2}\right) + \log\left(\frac{\mu_2}{\mu_1}\right).$$
(202)

It is also straightforward to check that the choice of  $w_1$  and  $w_2$  as in the statement of the theorem meets this upper bound and is thus optimal.

*Case ii*): In the second case where  $\chi_1 > \chi_2 > 1$ , we start as in *Case i*) and after optimization over  $\{\theta_{\bullet}\}$ , we can bound the sum-rate as

$$2E[R_{1}] + 2E[R_{2}] + 4\log(2) \\ \leq f\left(\frac{2\sqrt{\chi_{1}} \cdot \sin(\theta + \phi)}{(\chi_{1} - 1)(\sin^{2}(\theta) + \sin^{2}(\phi)) + 2}\right) + f\left(\frac{2\sqrt{\chi_{2}} \cdot \sin(\theta + \phi)}{(\chi_{2} - 1)(\sin^{2}(\theta) + \sin^{2}(\phi)) + 2}\right) \\ + 2\log\left(\frac{2\sqrt{\chi_{1}} \cdot \sin(\theta + \phi)}{(\chi_{1} - 1)\sin^{2}(\phi) + 1}\right) + 2\log\left(\frac{2\sqrt{\chi_{2}} \cdot \sin(\theta + \phi)}{(\chi_{2} - 1)\sin^{2}(\theta) + 1}\right)$$
(203)
$$= g\left(\frac{2\sqrt{\chi_{1}} \cdot \sin(\theta + \phi)}{(\chi_{1} - 1)(\sin^{2}(\theta) + \sin^{2}(\phi)) + 2}\right) + g\left(\frac{2\sqrt{\chi_{2}} \cdot \sin(\theta + \phi)}{(\chi_{2} - 1)(\sin^{2}(\theta) + \sin^{2}(\phi)) + 2}\right) \\ + 2\log\left(\frac{(\chi_{1} - 1)(\sin^{2}(\theta) + \sin^{2}(\phi)) + 2}{(\chi_{1} - 1)\sin^{2}(\phi) + 1}\right) + 2\log\left(\frac{(\chi_{2} - 1)(\sin^{2}(\theta) + \sin^{2}(\phi)) + 2}{(\chi_{2} - 1)\sin^{2}(\theta) + 1}\right).$$
(204)

The joint dependence between  $\theta$  and  $\phi$  in the right-hand side of (204) precludes the possibility of breaking down the double variable optimization of (204) into a pair of single variable optimizations. That is, the technique from *Case i*) fails here and this case needs to be studied differently.

The proof in this case follows in three steps. In the first step, when  $\chi_1 > \chi_2$ , we show that  $\mathcal{L}_2$  (defined as below) is maximized by  $\theta = \pi/2$  and  $\phi = 0$ :

$$\mathcal{L}_{2} \triangleq \frac{(\chi_{1} - 1)(\sin^{2}(\theta) + \sin^{2}(\phi)) + 2}{(\chi_{1} - 1)\sin^{2}(\phi) + 1} \cdot \frac{(\chi_{2} - 1)(\sin^{2}(\theta) + \sin^{2}(\phi)) + 2}{(\chi_{2} - 1)\sin^{2}(\theta) + 1}$$
(205)

$$=\frac{(\chi_1-1)(\chi_2-1)\left[\sin^2(\theta)+\sin^2(\phi)\right]^2+2(\chi_1+\chi_2-2)\left[\sin^2(\theta)+\sin^2(\phi)\right]+4}{\left[(\chi_1-1)\sin^2(\phi)+1\right]\left[(\chi_2-1)\sin^2(\theta)+1\right]}.$$
 (206)

For this, we set  $\sin^2(\theta) + \sin^2(\phi)$  to take a specific value of  $\alpha$ . Since the numerator is only a function of  $\alpha$ , maximizing  $\mathcal{L}_2$  is equivalent to minimizing the denominator of (206). There are two possible cases depending on whether  $\alpha \leq 1$  or  $1 < \alpha \leq 2$ . In the former case, it can be seen that the denominator is minimized by  $\phi = 0$  and  $\theta = \sin^{-1}(\sqrt{\alpha})$ , whereas in the latter case, it is minimized by  $\phi = \sin^{-1}(\sqrt{\alpha}-1)$  and  $\theta = \pi/2$ . Substituting these values, it can be seen that

$$\mathcal{L}_{2} \leq \begin{cases} \frac{(\chi_{1}-1)(\chi_{2}-1)\alpha^{2}+2\alpha(\chi_{1}+\chi_{2}-2)+4}{1+(\chi_{2}-1)\alpha} & \text{if } \alpha \leq 1\\ \frac{(\chi_{1}-1)(\chi_{2}-1)\alpha^{2}+2\alpha(\chi_{1}+\chi_{2}-2)+4}{\chi_{2}\cdot[1+(\chi_{1}-1)(\alpha-1)]} & \text{if } 1 < \alpha \leq 2. \end{cases}$$
(207)

A straightforward derivative calculation (using the critical fact that  $\chi_1 > \chi_2$ ) shows that while the right-hand side of (207) is increasing for  $\alpha \le 1$ , it is decreasing for  $1 < \alpha \le 2$ . In other words,

$$\mathcal{L}_2 \le \frac{(1+\chi_1) \cdot (1+\chi_2)}{\chi_2},\tag{208}$$

and this upper bound is achieved with  $\theta = \pi/2, \phi = 0$ .

In the second step, if  $\theta + \phi > \pi/2$ , we have

$$\frac{\sin(\theta + \phi)}{(\chi_i - 1)\left(\sin^2(\theta) + \sin^2(\phi)\right) + 2} \le \frac{\sin(\theta + \phi)}{\chi_i + 1} \le \frac{1}{\chi_i + 1}, \quad i = 1, 2$$
(209)

since  $\sin(\theta) > \sin(\pi/2 - \phi) = \cos(\phi)$ . Therefore,

 $2E[R_1] + 2E[R_2] + 4\log(2)$ 

$$\leq g\left(\frac{2\sqrt{\chi_1}\cdot\sin(\theta+\phi)}{(\chi_1-1)(\sin^2(\theta)+\sin^2(\phi))+2}\right) + g\left(\frac{2\sqrt{\chi_2}\cdot\sin(\theta+\phi)}{(\chi_2-1)(\sin^2(\theta)+\sin^2(\phi))+2}\right) + 2\log(\mathcal{L}_2)$$
(210)

$$\leq g\left(\frac{2\sqrt{\chi_1}}{\chi_1+1}\right) + g\left(\frac{2\sqrt{\chi_2}}{\chi_2+1}\right) + 2\log(\mathcal{L}_2),\tag{211}$$

$$2E[R_1] + 2E[R_2] \le f\left(\frac{2\sqrt{\chi_1}}{\chi_1 + 1}\right) + f\left(\frac{2\sqrt{\chi_2}}{\chi_2 + 1}\right) + \log\left(\frac{\chi_1}{\chi_2}\right),\tag{212}$$

where the second inequality follows from the monotonicity of  $g(\bullet)$  and the third from Step 1.

Note that (209) fails if  $\theta + \phi = \nu \le \pi/2$ . Thus, in the third step, we consider this possibility. Here, a straightforward manipulation shows that

$$\sin^{2}(\theta) + \sin^{2}(\nu - \theta) = \sin^{2}(\nu) - 2\sin(\theta)\sin(\nu - \theta)\cos(\nu) \le \sin^{2}(\nu),$$
 (213)

where the last inequality follows because  $\nu \leq \pi/2$ . Hence, we have

$$2E[R_1] + 2E[R_2] + 4\log(2) \le g\left(\frac{2\sqrt{\chi_1}\sin(\nu)}{(\chi_1 - 1)\sin^2(\nu) + 2}\right) + g\left(\frac{2\sqrt{\chi_2}\sin(\nu)}{(\chi_2 - 1)\sin^2(\nu) + 2}\right) + 2\log\left(\frac{[(\chi_1 - 1)\sin^2(\nu) + 2] \cdot [(\chi_2 - 1)\sin^2(\nu) + 2]}{[(\chi_1 - 1)\sin^2(\nu - \theta) + 1] \cdot [(\chi_2 - 1)\sin^2(\theta) + 1]}\right) \triangleq \mathcal{L}_3.$$
(214)

It can be easily seen that  $\mathcal{L}_3$  is maximized by  $\theta = \pi/2$  and  $\phi = 0$ . The upper bound is the same in both the cases  $\theta + \phi \le \pi/2$  and  $\theta + \phi > \pi/2$ . And this upper bound is met by the choice  $\theta = \pi/2$  and  $\phi = 0$  (that is,  $|\alpha| = 1$  and  $|\gamma| = 0$ ) and is hence optimal.

*Case iii*): Since the expression for the sum-rate is symmetric in  $\chi_1$  and  $\chi_2$ , an argument analogous to *Case ii*) completes the theorem in the case  $\chi_1 \leq \chi_2$ .

The optimal sum-rate in all the three cases is given by the unified expression

$$2E[R_1] + 2E[R_2] \xrightarrow{\rho \to \infty} f\left(\frac{2\sqrt{\chi_1}}{\chi_1 + 1}\right) + f\left(\frac{2\sqrt{\chi_2}}{\chi_2 + 1}\right) + \left|\log\left(\chi_1\right) - \log\left(\chi_2\right)\right|$$
(215)

where  $f(\bullet)$  is as defined in (126). This expression can be simplified as in the statement of the theorem.

## F. Comparison of Proof Techniques of Theorems 2 and 3

We first show that Theorem 2 reduces to Theorem 3 under the assumption that the eigenvectors of  $\Sigma_1$  and  $\Sigma_2$  coincide. For this, we set

$$\alpha' = \frac{\alpha\sqrt{\tau_1}}{X}, \quad \beta' = \frac{\beta\sqrt{\tau_2}}{X}, \quad \gamma' = \frac{\gamma\sqrt{\tau_1}}{Y}, \quad \delta' = \frac{\delta\sqrt{\tau_2}}{Y}$$
(216)

where X and Y are as in (152) and (153), respectively. Note that the above transformation is a bijection from the space  $\{\alpha, \beta, \gamma, \delta : |\alpha|^2 + |\beta|^2 = 1 = |\gamma|^2 + |\delta|^2\}$  to the space  $\{\alpha', \beta', \gamma', \delta' : |\alpha'|^2 + |\beta'|^2 = 1 = |\gamma'|^2 + |\delta'|^2\}$ . With this transformation, it can be checked that (156) reduces to (195). It can also be seen that the sum-rate expression in (45) reduces to that in (66) in both cases.

The technique pursued in the general case diverges from that in Appendix E in two ways. **Difference 1:** It can be easily checked that  $\tau_3 = 0$  if and only if the set of eigenvectors of  $\Sigma_1$ and  $\Sigma_2$  coincide. In general,  $\tau_3 \neq 0$  and  $\arg(\tau_3)$  could affect the sum-rate optimization. The first step in Appendix D is to show that this is not the case and  $\arg(\tau_3)$  plays no role in the optimization. This is done by exploiting the fact that the sum-rate optimization (see (14)) is a problem over  $\mathcal{G}(2, 1)$  and not over  $\operatorname{St}(2, 1)$ .

**Difference 2:** The second complication is that there is a definitive (and easily classifiable) comparative relationship between  $\tau_1 = \mathbf{v}_1^H \boldsymbol{\Sigma}_2^{-1} \mathbf{v}_1$  and  $\tau_2 = \mathbf{v}_2^H \boldsymbol{\Sigma}_2^{-1} \mathbf{v}_2$  in the special case. This comparative relationship does not generalize to the setting where the eigenvectors of  $\boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2$  are different.

Specifically, under Case i) of the discussion in Theorem 3,  $\tau_1 > \tau_2$  if and only if  $\chi_2 < 1$  whereas under Case ii),  $\tau_1 > \tau_2$  if and only if  $\chi_2 > 1$ . On the other hand, in the general case, all the three possibilities: i)  $\tau_1 > \tau_2$ , ii)  $\tau_1 < \tau_2$ , iii)  $\tau_1 = \tau_2$  can occur for appropriate choices of  $\Sigma_1$  and  $\Sigma_2$ . We now illustrate this with a numerical example. Let  $\Sigma_2$  be fixed such that

$$\Sigma_2 = \begin{bmatrix} 1 & -0.6897 \\ -0.6897 & 1 \end{bmatrix}.$$
 (217)

With the choice

$$\Sigma_1 = \begin{bmatrix} 1 & 0.8\\ 0.8 & 1 \end{bmatrix}, \tag{218}$$

it can be seen that  $\eta_1 = 5.8, \eta_2 = 0.1184, \tau_1 = 3.2222$  and  $\tau_2 = 0.5918$ , whereas if

$$\Sigma_1 = \begin{bmatrix} 1 & -0.8\\ -0.8 & 1 \end{bmatrix},\tag{219}$$

it can be seen that  $\eta_1 = 1.0653$ ,  $\eta_2 = 0.6444$ ,  $\tau_1 = 0.5918$  and  $\tau_2 = 3.2222$ . It can be seen that  $\eta_1 = 1.4603$ ,  $\eta_2 = 0.5397$  and  $\tau_1 = \tau_2 = 1.90725$  if  $\Sigma_1$  satisfies

$$\Sigma_1 = \begin{bmatrix} \frac{2}{3} & -0.34485\\ -0.34485 & \frac{1}{3} \end{bmatrix}.$$
 (220)

These differences imply that, in general, there exists no bijective transformation (as in (216)) to transform the objective function from the form in (156) to that in (195). Despite these issues, it would be of interest to pursue a theme that could unify the general case with the special case.

## G. Proof of Prop. 4

The proof in the low-SNR extreme is obvious. In the high-SNR extreme, we first note that the optimization problem over the choice of a pair  $(\mathbf{w}_1, \mathbf{w}_2)$  that results in a corresponding choice of  $(A_i, B_i, C_i)$  can be recast in the form of a two parameter optimization problem over  $(M_i, N_i)$  with  $M_i = \frac{A_i}{B_i}$  and  $N_i = \frac{C_i}{B_i}$  under the constraint that

$$0 \le N_i^2 \le M_i \le \chi_i = \frac{\lambda_1(\boldsymbol{\Sigma}_i)}{\lambda_2(\boldsymbol{\Sigma}_i)}.$$
(221)

For this, observe that for any given choice of  $(\mathbf{w}_1, \mathbf{w}_2)$ , the resultant  $(A_i, B_i, C_i)$  has to satisfy  $C_i^2 \leq A_i B_i$  (Cauchy-Schwarz inequality) and  $\frac{A_i}{B_i} \leq \chi_i$  (Ritz-Raleigh ratio) [47]. Thus, we have

$$\max_{\mathbf{w}_1, \mathbf{w}_2} \lim_{\rho \to \infty} E\left[R_i\right] \le \max_{0 \le N_i^2 \le M_i \le \chi_i} \lim_{\rho \to \infty} E\left[R_i\right].$$
(222)

Since the high-SNR expression for  $E[R_i]$  satisfies

$$2E[R_i] + 2\log(2) = g\left(\frac{2\sqrt{M_i - N_i^2}}{M_i + 1}\right) + 2\log(1 + M_i), \qquad (223)$$

and  $g(\bullet)$  is an increasing function, optimization over  $N_i$  which affects only the first term on the right-hand side implies that the optimal choice of  $N_i$  is zero. Plugging this choice and using the structure of  $g(\bullet)$  and  $f(\bullet)$  from (125) and (126) respectively, we have

$$E[R_i] \le \frac{M_i}{|M_i - 1|} \cdot |\log(M_i)|.$$
 (224)

The right-hand side of (224) is increasing in  $M_i$  since the derivative function satisfies

$$\frac{d\frac{M_i \cdot |\log(M_i)|}{|M_i - 1|}}{dM_i} = \begin{cases}
\frac{M_i - 1 - \log(M_i)}{(M_i - 1)^2} & \text{if } M_i > 1 \\
\frac{1}{2} & \text{if } M_i = 1 \\
\frac{\log\left(\frac{1}{M_i}\right) - (1 - M_i)}{(1 - M_i)^2} & \text{if } M_i < 1,
\end{cases}$$
(225)

where all the three pieces are positive and hence, the derivative function is smooth. The goal of maximizing  $M_i$  under the constraints of  $N_i = 0$  and unit-normedness of  $w_1$  and  $w_2$  is met by

Recall from (124) (the high-SNR expression) that  $E[R_i]$  is the sum of two terms. The increasing nature of  $g(\bullet)$  means that the first term of (124) is maximized when  $d_{\Sigma_i}(\mathbf{w}_1, \mathbf{w}_2)$  is maximized. That is, by the choice  $\{\mathbf{w}_1, \mathbf{w}_2\}$  as in (134)-(135). On the other hand, the second term as well as  $E[R_i]$  (which is the sum of the two terms) are maximized by the choice in (85). With this choice of beamforming vectors,  $d_{\Sigma_i}(\cdot, \cdot)$  can be written as

$$d_{\Sigma_i}(\mathbf{w}_{i,\mathsf{opt}},\mathbf{w}_{j,\mathsf{opt}}) = \frac{2\sqrt{\chi_i}}{\chi_i+1}.$$
(226)

Note that  $d_{\Sigma_i}(\cdot, \cdot)$  decreases as  $\chi_i$  increases.

#### REFERENCES

- G. J. Foschini, "Layered space-time architechture for wireless communication in a fading environment when using multielement antennas," *Bell Labs Tech. J.*, vol. 1, no. 2, pp. 41–59, 1996.
- [2] E. Visotsky and U. Madhow, "Space-time transmit precoding with imperfect feedback," *IEEE Trans. Inf. Theory*, vol. 47, no. 6, pp. 2632–2639, Sept. 2001.
- [3] A. J. Goldsmith, S. A. Jafar, N. Jindal, and S. Vishwanath, "Capacity limits of MIMO channels," *IEEE Journ. Sel. Areas in Commun.*, vol. 21, no. 5, pp. 684–702, June 2003.
- [4] D. Gesbert, H. Bolcskei, D. A. Gore, and A. J. Paulraj, "Outdoor MIMO wireless channels: Models and performance prediction," *IEEE Trans. Commun.*, vol. 50, no. 12, pp. 1926–1934, Dec. 2002.
- [5] V. V. Veeravalli, Y. Liang, and A. M. Sayeed, "Correlated MIMO Rayleigh fading channels: capacity, optimal signaling and asymptotics," *IEEE Trans. Inf. Theory*, vol. 51, no. 6, pp. 2058–2072, June 2005.
- [6] Q. H. Spencer, C. B. Peel, A. L. Swindlehurst, and M. Haardt, "An introduction to the multi-user MIMO downlink," *IEEE Commun. Magaz.*, vol. 42, no. 10, pp. 60–67, Oct. 2004.
- [7] D. Gesbert, M. Kountouris, R. W. Heath, Jr., C.-B. Chae, and T. Salzer, "Shifting the MIMO paradigm: from single user to multiuser communications," *IEEE Sig. Proc. Magaz.*, vol. 24, no. 5, pp. 36–46, Oct. 2007.
- [8] B. Hassibi and M. Sharif, "Fundamental limits in MIMO broadcast channels," *IEEE Journ. Sel. Areas in Commun.*, vol. 25, no. 7, pp. 1333–1344, Sept. 2007.
- [9] G. Caire, N. Jindal, M. Kobayashi, and N. Ravindran, "Multiuser MIMO achievable rates with downlink training and channel state feedback," *IEEE Trans. Inf. Theory*, vol. 56, no. 6, pp. 2845–2866, June 2010.
- [10] G. Caire and S. Shamai, "On the achievable throughput of a multiantenna Gaussian broadcast channel," *IEEE Trans. Inf. Theory*, vol. 49, no. 7, pp. 1691–1706, July 2003.
- [11] P. Viswanath and D. N. C. Tse, "Sum capacity of the vector Gaussian broadcast channel and downlink-uplink duality," *IEEE Trans. Inf. Theory*, vol. 49, no. 8, pp. 1912–1921, Aug. 2003.
- [12] S. Vishwanath, N. Jindal, and A. J. Goldsmith, "Duality, achievable rates and sum rate capacity of Gaussian MIMO broadcast channel," *IEEE Trans. Inf. Theory*, vol. 49, no. 10, pp. 2658–2668, Oct. 2003.
- [13] N. Jindal, S. Vishwanath, and A. J. Goldsmith, "On the duality of Gaussian multiple-access and broadcast channels," *IEEE Trans. Inf. Theory*, vol. 50, no. 5, pp. 768–783, May 2004.
- [14] W. Yu and J. M. Cioffi, "Sum capacity of Gaussian vector broadcast channels," *IEEE Trans. Inf. Theory*, vol. 50, no. 9, pp. 1875–1892, Sept. 2004.

- [15] M. Schubert and H. Boche, "Solution of multiuser downlink beamforming problem with individual SINR constraint," *IEEE Trans. Veh. Tech.*, vol. 53, no. 1, pp. 18–28, Jan. 2004.
- [16] H. Weingarten, Y. Steinberg, and S. Shamai, "The capacity region of the Gaussian multiple-input multiple-output broadcast channel," *IEEE Trans. Inf. Theory*, vol. 52, no. 9, pp. 3936–3964, Sept. 2006.
- [17] M. H. M. Costa, "Writing on dirty paper," IEEE Trans. Inf. Theory, vol. 29, no. 3, pp. 439-441, May 1983.
- [18] A. Bennatan, D. Burshtein, G. Caire, and S. Shamai, "Superposition coding for side-information channels," *IEEE Trans. Inf. Theory*, vol. 52, no. 5, pp. 1872–1889, May 2006.
- [19] R. F. H. Fischer, C. Windpassinger, A. Lampe, and J. B. Huber, "MIMO precoding for decentralized receivers," *Proc. IEEE Intern. Symp. Inf. Theory*, p. 496, 2002.
- [20] M. Joham, W. Utschick, and J. Nossek, "Linear transmit processing in MIMO communications systems," *IEEE Trans. Sig. Proc.*, vol. 53, no. 8, pp. 2700–2712, Aug. 2005.
- [21] F. Boccardi, F. Tosato, and G. Caire, "Precoding Schemes for the MIMO-GBC," Proc. 2006 Int. Zurich Seminar on Commun., pp. 10–13, Feb. 2006.
- [22] B. M. Hochwald, C. B. Peel, and A. L. Swindlehurst, "A vector-perturbation technique for near-capacity multiantenna multiuser communication - Part II: perturbation," *IEEE Trans. Commun.*, vol. 53, no. 3, pp. 537–544, Mar. 2005.
- [23] T. Yoo and A. J. Goldsmith, "On the optimality of multiantenna broadcast scheduling using zero-forcing beamforming," *IEEE Journ. Sel. Areas in Commun.*, vol. 24, no. 3, pp. 528–541, Mar. 2006.
- [24] Q. H. Spencer, A. L. Swindlehurst, and M. Haardt, "Zero-forcing methods for downlink spatial multiplexing in multi-user MIMO channels," *IEEE Trans. Sig. Proc.*, vol. 52, no. 2, pp. 461–471, Feb. 2004.
- [25] A. Wiesel, Y. C. Eldar, and S. Shamai, "Zero forcing precoding and generalized inverses," *IEEE Trans. Sig. Proc.*, vol. 56, no. 9, pp. 4409–4418, Sept. 2008.
- [26] C.-B. Chae, D. Mazzarese, N. Jindal, and R. W. Heath, Jr., "Coordinated beamforming with limited feedback in the MIMO broadcast channel," *IEEE Journ. Sel. Areas in Commun.*, vol. 26, no. 8, pp. 1505–1515, Oct. 2008.
- [27] C.-B. Chae, S. Shim, and R. W. Heath, Jr., "Block diagonalized vector perturbation for multi-user MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 7, no. 11, pp. 4051–4057, Nov. 2008.
- [28] V. Raghavan, J. H. Kotecha, and A. M. Sayeed, "Why does the Kronecker model result in misleading capacity estimates?," *IEEE Trans. Inf. Theory*, vol. 56, no. 10, pp. 4843–4864, Oct. 2010.
- [29] T. Y. Al-Naffouri, M. Sharif, and B. Hassibi, "How much does transmit correlation affect the sum-rate scaling of MIMO Gaussian broadcast channels?," *IEEE Trans. Commun.*, vol. 57, no. 2, pp. 562–572, Feb. 2009.
- [30] M. Kountouris, D. Gesbert, and L. Pittman, "Transmit correlation-aided opportunistic beamforming and scheduling," Proc. European Sig. Proc. Conf., pp. 2409–2414, Sept. 2006.
- [31] D. Hammarwall, M. Bengtsson, and B. E. Ottersten, "Acquiring partial CSI for spatially selective transmission by instantaneous channel norm feedback," *IEEE Trans. Sig. Proc.*, vol. 56, no. 3, pp. 1188–1204, Mar. 2008.
- [32] D. Hammarwall, M. Bengtsson, and B. E. Ottersten, "Utilizing the spatial information provided by channel norm feedback in SDMA systems," *IEEE Trans. Sig. Proc.*, vol. 56, no. 7-2, pp. 3278–3293, July 2008.
- [33] M. Trivellato, F. Boccardi, and H. Huang, "On transceiver design and channel quantization for downlink multiuser MIMO systems with limited feedback," *IEEE Journ. Sel. Areas in Commun.*, vol. 6, no. 8, pp. 1494–1504, Oct. 2008.
- [34] V. Raghavan and S. V. Hanly, "Limited feedback codebook design for interference management," Presented at the Inf. Theory and Appl. Workshop, San Diego, Feb. 2010.
- [35] B. Clerckx, G. Kim, and S. Kim, "Correlated fading in broadcast MIMO channels: curse or blessing?," Proc. IEEE Global Telecommun. Conf., pp. 3830–3834, Dec. 2008.

- [36] R. de Francisco, C. Simon, D. T. M. Slock, and G. Leus, "Beamforming for correlated broadcast channels with quantized channel state information," *Proc. IEEE Workshop Sig. Proc. Adv. in Wireless Commun., Brazil*, pp. 161–165, July 2008.
- [37] A. Wiesel, Y. C. Eldar, and S. Shamai, "Linear precoding via conic optimization for fixed MIMO receivers," *IEEE Trans. Sig. Proc.*, vol. 54, no. 1, pp. 161–176, Jan. 2006.
- [38] V. S. Annapureddy and V. V. Veeravalli, "Sum capacity of MIMO interference channels in the low interference regime," *To appear in the IEEE Trans. Inf. Theory*, 2011.
- [39] V. Raghavan and S. V. Hanly, "Statistical beamformer design for the two-antenna interference channel," *Proc. IEEE Intern. Symp. Inf. Theory*, pp. 2278–2282, June 2010.
- [40] M. Kobayashi and G. Caire, "An iterative water-filling algorithm for maximum weighted sum-rate of Gaussian MIMO-BC," *IEEE Journ. Sel. Areas in Commun.*, vol. 24, no. 8, pp. 1640–1646, Aug. 2006.
- [41] N. Merhav, G. Kaplan, A. Lapidoth, and S. Shamai, "On information rates for mismatched decoders," *IEEE Trans. Inf. Theory*, vol. 40, no. 6, pp. 1953–1967, June 1994.
- [42] A. Edelman, T. Arias, and S. Smith, "The geometry of algorithms with orthogonality constraints," SIAM J. Matrix Anal. Appl., vol. 20, no. 2, pp. 303–353, Apr. 1999.
- [43] B. Hassibi and T. L. Marzetta, "Multiple-antennas and isotropically random unitary inputs: the received signal density in closed form," *IEEE Trans. Inf. Theory*, vol. 48, no. 6, pp. 1473–1484, June 2002.
- [44] M. Abramowitz and I. A. Stegun, Handbook of Mathematical Functions with Formulas, Graphs and Mathematical Tables, National Bureau of Standards, USA, 10th edition, 1972.
- [45] I. S. Gradshteyn and I. M. Ryzhik, Table of Integrals, Series, and Products, Academic Press, NY, 4th edition, 1965.
- [46] G. W. Stewart, Matrix Algorithms: Eigensystems, SIAM Publishers, USA, 2001.
- [47] C. R. Rao and M. B. Rao, *Matrix Algebra and its Applications to Statistics and Econometrics*, World Scientific Publishing Co. Pte. Ltd., Singapore, 1998.
- [48] C. A. Floudas, Nonlinear and Mixed-Integer Optimization: Fundamentals and Applications, Oxford University Press, UK, 1995.
- [49] V. Raghavan, V. V. Veeravalli, and R. W. Heath, Jr., "Reduced rank signaling in spatially correlated MIMO channels," *Proc. IEEE Intern. Symp. Inf. Theory*, pp. 1081–1085, July 2007.
- [50] V. Raghavan, A. M. Sayeed, and V. V. Veeravalli, "Semiunitary precoding for spatially correlated MIMO channels," *IEEE Trans. Inf. Theory*, vol. 57, no. 3, pp. 1284–1298, Mar. 2011.
- [51] K. K. Mukkavilli, A. Sabharwal, E. Erkip, and B. Aazhang, "On beamforming with finite rate feedback in multiple antenna systems," *IEEE Trans. Inf. Theory*, vol. 49, no. 10, pp. 2562–2579, Oct. 2003.