# Efficient Weighted Sum Rate Maximization With Linear Precoding

Christian Guthy, *Student Member, IEEE*, Wolfgang Utschick, *Senior Member, IEEE*,
Raphael Hunger, *Student Member, IEEE*, and Michael Joham, *Member, IEEE*

*Abstract*—Achieving the boundary of the capacity region in the multiple-input multiple-output (MIMO) broadcast channel requires the use of dirty paper coding (DPC). As practical nearly optimum implementations of DPC are computationally complex, purely linear approaches are often used instead. However, in this case, the problem of maximizing a weighted sum rate constitutes a nonconvex and, in most cases, also a combinatorial optimization problem. In this paper, we present two heuristic nearly optimum algorithms with reduced computational complexity. For this purpose, a lower bound for the weighted sum rate under linear zero-forcing constraints is used. Based on this bound, both greedy algorithms successively allocate data streams to users. In each step, the user is determined that is given an additional data stream such that the increase in weighted sum rate becomes maximum. Thereby, the data stream allocations and filters obtained in the previous steps are kept fixed and only the filter corresponding to the additional data stream is optimized. The first algorithm determines the receive and transmit filters directly in the downlink. The other algorithm operates in the dual uplink, from which the downlink transmit and receive filters can be obtained via the general rate duality leading to nonzero-forcing in the downlink. Simulation results reveal marginal performance losses compared to more complex algorithms.

*Index Terms*—Broadcast channel, linear precoding, multiple-input multiple-output (MIMO) systems.

## I. INTRODUCTION

A popular method to achieve a point on the boundary of the capacity region of the multiple-input multiple-output (MIMO) broadcast channel, which has recently been found in [1], is to solve a weighted sum rate maximization. Viswanathan *et al.* have first presented a solution to this problem for the MIMO broadcast channel [2] and more efficient algorithms have been developed afterwards in [3], [4] and in [5] for multiple-input single-output (MISO) systems. By varying the weights, different points on the boundary can be achieved. In other applications, where a weighted sum rate maximization occurs, the weights correspond to the priorities of the users or the length

of the queues of the users [5]. Other quality-of-service (QoS) constrained problems in the MIMO broadcast channel are often solved by an iterative application of the weighted sum rate maximization, i.e., the appropriate weights are determined fulfilling the given QoS constraints and optimizing the objective function in an iterative manner. That applies for the problem of maximizing sum rate under a total power and relative rate constraints [6], the maximization of a weighted sum rate under a power constraint and minimum rate requirements [7], and the minimization of transmit power required to satisfy minimum rate requirements [8]. All these algorithms operate on the boundary of the capacity region and therefore require the use of dirty paper coding (DPC) [9]. To achieve the capacity region with DPC in practice, numerically complex methods such as vector precoding [10] or the coding scheme from [11] must be employed. The latter method also exhibits long encoding delays. When some losses in sum rate compared to the optimum are acceptable, the complexity of the implementation of DPC can be reduced, for example through the use of Tomlinson-Harashima precoding (THP) [12], [13]. The reasons for the suboptimality of this scheme are for example explained in [14]. Nevertheless, THP still exhibits practical challenges, such as the implementation of modulo operators at all receivers due to the dynamics of the received signals. For this reason, linear approaches, which minimize or cancel completely the interference between data streams by linear signal processing and therefore avoid the use of DPC, are of high practical relevance. However, in contrast to the DPC case, determining the optimum transmit and receive filters requires the solution of a nonconvex optimization problem. In case the sum of receive antennas at the terminals in the system is larger than the number of transmit antennas at the base station—a setup that is very likely in a practical system—the optimization problem additionally becomes combinatorial. Obviously that is the case when zero-forcing is employed at the transmitter, as the number of data streams that can be multiplexed in space is limited by the number of transmit antennas and a combinatorial search for the optimum allocation of data streams to users is required. In [15], it has been shown that such a combinatorial search is also required for the optimum linear solution without zero-forcing constraints, as the (weighted) sum rate utility is nonconcave. Since globally optimum solutions have prohibitive complexity, one has to go for locally optimum ones which are prone to converge to local optima with the same stream configuration they have been initialized with, see the precoder based projected gradient approach in [16]. Hence, all possible stream configurations would in principle have to be probed to come to an almost globally optimum

performance which then leads to the combinatorial search. Furthermore, the optimum rate region achievable with linear precoding still constitutes an open problem. Only for the special case of two users and single antenna receivers a solution to achieve points on the boundary of the rate region has been presented in [17]. The problem of weighted sum rate maximization with linear precoding is claimed to be solved by the algorithm of [16]. However, convergence to the optimum solution is not guaranteed and strongly depends on the initialization. The algorithm works iteratively, requires formulating and solving a geometric program in each step and is based on a repeated transformation from the dual multiple access channel to the broadcast channel and back enabled by the single data stream duality of [18]. A general concept for the duality between these channels with joint coding has been developed in [19]. Due to the involved repeated formulations and solutions of geometric programs, the algorithm of [16] exhibits a considerable amount of computational complexity. The approach presented in [20] also requires an iterative application of geometric programming. Besides aiming at maximizing a weighted sum rate, [20] considers the problem of feedback reduction. A projected gradient method is used to solve the weighted sum rate maximization problem with linear precoding in [4] and [15].

To avoid the complexity associated with the power allocation, which requires the use of geometric programming in [16] and [20], linear zero-forcing techniques will be employed in this paper. Interference is thereby suppressed completely by linear signal processing. We first apply linear zero-forcing directly in the downlink and then in the dual uplink, where from the filters in the downlink can be computed via the general duality of [19]. In case the number of transmit antennas is larger than or equal to the total number of receive antennas in the system, block diagonalization (BD) [21] can be applied to determine the precoders in the downlink for complete interference suppression, that reduces to zero-forcing beamforming (ZFBF) [22] in MISO systems. However, if the number of transmit antennas is smaller than the total number of receive antennas, not all users get as many data streams as they have receive antennas. Thus, a search for the best allocation of data streams to users is inevitable. Optimally this would imply an exhaustive search over all possible allocations. As such a search becomes infeasible already for a moderate number of users, several heuristic methods for this problem have been proposed.

For the special case of equal weights for all users, i.e., maximization of sum rate, heuristic user selection methods for BD have been proposed in [23], [24], and [25]. For ZFBF in MISO systems, this selection is performed in a greedy manner in [26], i.e., beginning with the strongest user, a data stream is given to the user that leads to the strongest increase in sum rate in each step. A low complexity implementation of the algorithm is presented in [27]. In [28], the required search for the best user in each step is simplified by excluding users due to their spatial channel properties. These approaches for MISO systems can be extended to MIMO by simply applying the left singular vectors of each user's channel matrix as receive filters and considering each resulting product of singular value and corresponding right singular vector as a virtual user channel. A more advanced extension of the algorithm from [26] to MIMO can be found in

[29], where the receive filters are initialized with the left singular vectors and adjusted in case when more than one data stream is allocated to one user. An algorithm which includes the receive filters into the successive optimization and is additionally less complex than the previous approaches at no performance loss is described in [30] and [31]. The concept of successive filter determination is also utilized in [32] for sum rate maximization in the downlink as well as the uplink. As the algorithms in [32] rely on an *a priori* fixed power allocation, their application to the weighted sum rate maximization considered in this paper is not straightforward. In [33], the concept of successive allocation is extended to the problem of weighted sum rate maximization for MISO systems, where the weighted sum rate is approximated by a lower bound to avoid computing sum rates explicitly for each tested user in every step. As in [33] for MISO systems, we also employ a successive allocation of data streams to users, where in each step not only the user but also the corresponding filter at the user that lead to the strongest increase in weighted sum rate are determined. Since even with this simplification the resulting optimization problem is still too complex, we use a lower bound for the weighted sum rate. By maximizing this lower bound, the resulting optimization can be solved via the computations of generalized eigenvectors and in contrast to state-of-the-art approaches, no iterative application of the algorithm is required at marginal performance loss. The proposed allocation can be applied directly in the downlink as well as in the dual uplink.

The outline of the paper is as follows: After explaining the system model in Section II, the main concept of the proposed linear zero-forcing scheme for weighted sum rate maximization is explained in Section III. The algorithm operating directly in the downlink is presented at a glance in Section IV and the algorithm based on the duality between uplink and downlink is described in Section V. Simulation results are shown in Section VI and the paper is concluded in Section VII.

*Notation:* Bold lower and uppercase letters denote vectors and matrices, respectively. $(\bullet)^{\mathrm{T}}$ and $(\bullet)^{\mathrm{H}}$ describe the transpose and the Hermitian of a matrix, respectively. $\lambda_{\max}(\boldsymbol{A})$, $\mathrm{tr}(\boldsymbol{A})$, $|\boldsymbol{A}|$, $\|\boldsymbol{A}\|_{\mathrm{F}}$, and $(\boldsymbol{A})_{i,j}$ are the maximum eigenvalue, the trace, the determinant, the Frobenius norm, and the element in row $i$ and column $j$ of the matrix $\boldsymbol{A}$, respectively. $\boldsymbol{A}^{+}$ denotes the Moore-Penrose pseudoinverse of the matrix $\boldsymbol{A}$, $\mathrm{vec}(\boldsymbol{A})$ stacks the columns of the matrix $\boldsymbol{A}$ in one vector and $\mathrm{diag}(a_1, \ldots, a_i)$ denotes a diagonal matrix with the elements $a_1, \ldots, a_i$ on its diagonal. $\boldsymbol{I}_i$ is the $i \times i$ identity matrix, $\boldsymbol{0}_{i,j}$ is the $i \times j$ zero matrix, and $\boldsymbol{e}_j$ denotes the $j$th canonical unit vector. $\mathrm{null}\{\boldsymbol{A}\}$ and $(\mathrm{null}\{\boldsymbol{A}\})^{\perp}$ denote the nullspace of the matrix $\boldsymbol{A}$ and the orthogonal complement to this nullspace, respectively.

## II. SYSTEM MODEL

We consider a multiuser MIMO system with one base station and $K$ noncooperative users. The base station is equipped with $N$ antennas and the number of antennas at user $k$ is denoted by $r_k$. The system model of the broadcast channel is depicted in Fig. 1. $\boldsymbol{P}_k \in \mathbb{C}^{N \times d_k}$ denotes the $k$th user's precoding matrix in the broadcast channel, where $d_k \leq r_k$ accounts for the number of data streams of user $k$. $\boldsymbol{H}_k \in \mathbb{C}^{r_k \times N}$ and $\boldsymbol{B}_k^{\mathrm{H}} \in \mathbb{C}^{d_k \times r_k}$ are the channel and the receive filter of user $k$, respectively. We assume that each user has perfect knowledge of its channel
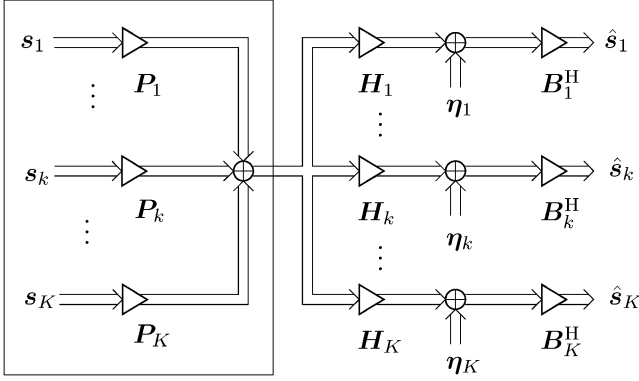
Fig. 1. System model of the MIMO broadcast channel.

matrix $\boldsymbol{H}_k$ and all users' channels are perfectly known at the base station. $\boldsymbol{\eta}_k \in \mathbb{C}^{r_k}$ denotes the additive noise experienced by user $k$, where we assume in the following that $\boldsymbol{\eta}_k$ is circularly symmetric Gaussian with zero mean and identity covariance matrix, i.e., $\boldsymbol{\eta}_k \sim \mathcal{CN}(\boldsymbol{0}, \boldsymbol{I}_{r_k})$. All noise signals are uncorrelated with all input signals $\boldsymbol{s}_1, \ldots, \boldsymbol{s}_K$, which have been encoded with Gaussian codebooks, i.e., $\boldsymbol{s}_k \sim \mathcal{CN}(\boldsymbol{0}, \boldsymbol{I}_{d_k})$. The received signal $\hat{\boldsymbol{s}}_k$ of user $k$ in the downlink is then given by

$$\hat{\boldsymbol{s}}_k = \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_k \boldsymbol{s}_k + \sum_{\substack{j=1 \\ j \neq k}}^{K} \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_j \boldsymbol{s}_j + \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{\eta}_k.$$

When no DPC is used, the $k$th user's rate computes according to

$$R_{k,\mathrm{DL}} = \log_2 \frac{\left| \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{B}_k + \sum_{j=1}^{K} \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_j \boldsymbol{P}_j^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \boldsymbol{B}_k \right|}{\left| \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{B}_k + \sum_{\substack{j=1 \\ j \neq k}}^{K} \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_j \boldsymbol{P}_j^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \boldsymbol{B}_k \right|}. \quad (1)$$

The average power constraint at the base station implies that $\sum_{k=1}^{K} \mathrm{tr}(\boldsymbol{P}_k \boldsymbol{P}_k^{\mathrm{H}}) \leq P_{\mathrm{Tx}}$.

## III. SUCCESSIVE WEIGHTED SUM RATE MAXIMIZATION

We consider the maximization of a weighted sum of the users' rates employing linear signal processing only, which leads to the following optimization problem for the broadcast channel:

$$\underset{\{\boldsymbol{P}_k, \boldsymbol{B}_k\}_{k=1,\ldots,K}}{\text{maximize}} \sum_{k=1}^{K} \mu_k R_{k,\mathrm{DL}} = \underset{\{\boldsymbol{P}_k, \boldsymbol{B}_k\}_{k=1,\ldots,K}}{\text{maximize}} \sum_{k=1}^{K} \mu_k$$
$$\times \log_2 \frac{\left| \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{B}_k + \sum_{j=1}^{K} \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_j \boldsymbol{P}_j^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \boldsymbol{B}_k \right|}{\left| \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{B}_k + \sum_{\substack{j=1 \\ j \neq k}}^{K} \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_j \boldsymbol{P}_j^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \boldsymbol{B}_k \right|}$$
$$\text{s.t. } \sum_{k=1}^{K} \mathrm{tr}\left(\boldsymbol{P}_k \boldsymbol{P}_k^{\mathrm{H}}\right) \leq P_{\mathrm{Tx}} \quad (2)$$

with *a priori* given weights $\mu_1, \ldots, \mu_K$. Note that in (2) $d_k = r_k$ for all users $k \in \{1, \ldots, K\}$. The case that a user can transmit over less than $r_k$ active data streams is modeled by allocating zero power to the inactive data streams. Equation (2) constitutes a nonconvex optimization problem even in the dual uplink,

for which to the authors' best knowledge no optimum solution has been found so far. We, therefore, propose to simplify the problem in three steps.

*1) Introduction of Zero-Forcing Constraints:* We impose zero-forcing constraints such that the data streams of different users do not interfere with each other. Although at low SNR zero-forcing suffers from noise enhancement, zero-forcing is optimum when interstream interference becomes dominant at high SNR. Introducing these zero-forcing constraints, Optimization (2) reduces to

$$\underset{\{\boldsymbol{P}_k, \boldsymbol{B}_k\}_{k=1,\ldots,K}}{\text{maximize}} \sum_{k=1}^{K} \mu_k$$
$$\times \log_2 \left| \boldsymbol{I}_{d_k} + \left( \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{B}_k \right)^{-1} \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_k \boldsymbol{P}_k^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \boldsymbol{B}_k \right|$$
$$\text{s.t. } \sum_{k=1}^{K} \mathrm{tr}\left(\boldsymbol{P}_k \boldsymbol{P}_k^{\mathrm{H}}\right) \leq P_{\mathrm{Tx}}$$
$$\boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_j = \boldsymbol{0}_{d_k, d_j}, \ \forall k, \forall j \neq k. \quad (3)$$

The zero-forcing constraints in (3) imply that at most $N$ data streams can be active, i.e., obtain nonzero power. Consequently, an exhaustive search for the optimum allocation of data streams to users would be necessary if the global optimum shall be found. Considering the fact that there are $\sum_{j=1}^{\min\{N,r\}} \binom{r}{j}$ possible allocations, where $r = \sum_{k=1}^{K} r_k$, this problem becomes infeasible already with a moderate number of users. Furthermore the problem of determining the optimum transmit and receive filters for a fixed allocation is still nonconvex.

In (3) $\boldsymbol{B}_k$ can be right-hand side (RHS) multiplied with any invertible matrix $\boldsymbol{C}_k \in \mathbb{C}^{d_k \times d_k}$ without changing the objective function and the constraints. Consequently, if $\boldsymbol{B}_k$ maximizes (3), so does $\tilde{\boldsymbol{B}}_k = \boldsymbol{B}_k \boldsymbol{C}_k$. We can therefore restrict $\boldsymbol{B}_k$ to have orthonormal columns without changing the optimization problem to simplify the rate expression in (3). Thus, (3) is equivalent to

$$\underset{\{\boldsymbol{P}_k, \boldsymbol{B}_k\}_{k=1,\ldots,K}}{\text{maximize}} \sum_{k=1}^{K} \mu_k \log_2 \left| \boldsymbol{I}_{d_k} + \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_k \boldsymbol{P}_k^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \boldsymbol{B}_k \right|$$
$$\text{s.t. } \sum_{k=1}^{K} \mathrm{tr}\left(\boldsymbol{P}_k \boldsymbol{P}_k^{\mathrm{H}}\right) \leq P_{\mathrm{Tx}}$$
$$\boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_j = \boldsymbol{0}_{d_k, d_j}, \ \forall k, \forall j \neq k, \quad \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{B}_k = \boldsymbol{I}_{d_k}, \forall k. \quad (4)$$

Still, the solution of (4) is not unique with respect to $\boldsymbol{B}_k$ and $\boldsymbol{P}_k$ as multiplying these matrices RHS with any $d_k \times d_k$ orthonormal matrix does not change the objective function nor the constraints. We can therefore restrict the matrices $\boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_k$ to be diagonal, which leads to the following optimization problem

$$\underset{\{\boldsymbol{P}_k, \boldsymbol{B}_k\}_{k=1,\ldots,K}}{\text{maximize}} \sum_{k=1}^{K} \mu_k \log_2 \left| \boldsymbol{I}_{d_k} + \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_k \boldsymbol{P}_k^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \boldsymbol{B}_k \right|$$
$$\text{s.t. } \sum_{k=1}^{K} \mathrm{tr}\left(\boldsymbol{P}_k \boldsymbol{P}_k^{\mathrm{H}}\right) \leq P_{\mathrm{Tx}}, \ \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_j = \boldsymbol{0}_{d_k, d_j}, \forall k, \forall j \neq k$$
$$\boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{B}_k = \boldsymbol{I}_{d_k}, \forall k, \quad \boldsymbol{B}_k^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{P}_k \text{ diagonal}, \forall k. \quad (5)$$

This way, the power allocation over the resulting interference-free data streams is drastically simplified as it can be found with weighted water-filling and all zero-forcing constraints in (5) can be incorporated easily into the objective function as described in the following. Consider the case that $i$ data streams are active, which are indexed by $j = 1, \ldots, i$. $\boldsymbol{b}_j^{\mathrm{H}}$ denotes the receive filter for the $j$th data stream and $\pi(j)$ is the user to which the $j$th data stream is allocated to. Note that $\boldsymbol{b}_j$ therefore corresponds to one column of $\boldsymbol{B}_{\pi(j)}$ in our system model. Then the composite channel matrix $\boldsymbol{H}_{\mathrm{co},i}$ can be defined as

$$\boldsymbol{H}_{\mathrm{co},i} = \begin{bmatrix} \boldsymbol{b}_1^{\mathrm{H}} \boldsymbol{H}_{\pi(1)} \\ \vdots \\ \boldsymbol{b}_i^{\mathrm{H}} \boldsymbol{H}_{\pi(i)} \end{bmatrix} \in \mathbb{C}^{i \times N}. \tag{6}$$

From now on, $d_k$ corresponds to the number of active data streams, i.e. the columns of $\boldsymbol{B}_k$ and $\boldsymbol{P}_k$ corresponding to inactive data streams are omitted. Furthermore we subsume the precoders of the active data streams into the effective precoder $\boldsymbol{P}_{\mathrm{eff},i}$, from which each precoding matrix $\boldsymbol{P}_k$ can be obtained by extracting all columns that correspond to data streams allocated to user $k$[1]. The effective precoder $\boldsymbol{P}_{\mathrm{eff},i}$ can be derived from the composite channel matrix with the zero-forcing constraints according to

$$\boldsymbol{P}_{\mathrm{eff},i} = \boldsymbol{H}_{\mathrm{co},i}^{+} \boldsymbol{\Lambda}_i \boldsymbol{\Gamma}_i^{\frac{1}{2}}. \tag{7}$$

The pseudoinverse $\boldsymbol{H}_{\mathrm{co},i}^{+}$ is used to diagonalize the composite channel. This precoder structure leads to the optimum solution of (5) (e.g., [34, Theorem 1]). The diagonal matrix $\boldsymbol{\Lambda}_i = \mathrm{diag}(\lambda_{i,1}, \ldots, \lambda_{i,i})$ is necessary to set the column norms of the pseudoinverse multiplied by $\boldsymbol{\Lambda}_i$ to unity, i.e.,

$$\lambda_{i,j} = \frac{1}{\left\| \boldsymbol{H}_{\mathrm{co},i}^{+} \boldsymbol{e}_j \right\|_2} \tag{8}$$

such that the power can be allocated to the subchannels according to $\boldsymbol{\Gamma}_i = \mathrm{diag}(\gamma_{i,1}, \ldots, \gamma_{i,i})$. The weighted sum rate after step $i$ can be written as [cf. (5)]

$$R_{\mathrm{wsr},i}\left(\pi(1), \ldots, \pi(i), \boldsymbol{b}_1, \ldots, \boldsymbol{b}_i\right)$$
$$= \sum_{j=1}^{i} \mu_{\pi(j)} \log_2 \left(1 + \gamma_{i,j} \lambda_{i,j}^2\right) \tag{9}$$

where the $\lambda_{i,j}$'s depend on the set of allocated users $\bigcup_{j=1}^{i}\{\pi(j)\}$ and on the receive filters $\boldsymbol{b}_1^{\mathrm{H}}, \ldots, \boldsymbol{b}_i^{\mathrm{H}}$ as shown in (8) and (6). Maximizing (9) with respect to the nonnegative $\gamma_{i,j}$'s under a total transmit power constraint leads to a concave optimization problem. By solving the Karush-Kuhn-Tucker conditions [35], we obtain

$$\gamma_{i,j} = \max \left\{ \mu_{\pi(j)} \eta_i - \frac{1}{\lambda_{i,j}^2}, 0 \right\} \tag{10}$$

where $\eta_i$ is chosen such that the transmit power constraint is fulfilled with equality, i.e., $\sum_{j=1}^{i} \gamma_{i,j} = P_{\mathrm{Tx}}$.

[1]All columns $j = 1, \ldots, i$ of $\boldsymbol{P}_{\mathrm{eff},i}$ for which $\pi(j) = k$ are therefore collected in $\boldsymbol{P}_k$.

*2) Successive Data Stream Allocation and Filter Computation:* To avoid the full combinatorial search necessary to solve (5), we follow a greedy successive allocation scheme along the lines of [26] and [33], i.e., the user allocation of the previous steps is kept fixed. The resulting problem in step $i$ can therefore be written as

$$\{\pi(i), \boldsymbol{b}_1, \ldots, \boldsymbol{b}_i\}$$
$$= \underset{\{k, \hat{\boldsymbol{b}}_1, \ldots, \hat{\boldsymbol{b}}_i\}}{\arg \max} R_{\mathrm{wsr},i}\left(\pi(1), \ldots, \pi(i-1), k, \hat{\boldsymbol{b}}_1, \ldots, \hat{\boldsymbol{b}}_i\right)$$
$$\text{s.t.} \ \hat{\boldsymbol{b}}_j^{\mathrm{H}} \hat{\boldsymbol{b}}_j = 1, \forall j \in \{1, \ldots, i\}$$
$$\hat{\boldsymbol{b}}_j^{\mathrm{H}} \hat{\boldsymbol{b}}_\ell = 0, \forall \ell \neq j \text{ with } \pi(\ell) = \pi(j), \forall j. \tag{11}$$

For $i = 1$, (11) reduces to

$$\{\pi(1), \boldsymbol{b}_1\} = \underset{\{k, \hat{\boldsymbol{b}}_1\}}{\arg \max} \mu_k \log_2 \left(1 + P_{\mathrm{Tx}} \hat{\boldsymbol{b}}_1^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{H}_k^{\mathrm{H}} \hat{\boldsymbol{b}}_1\right)$$
$$\text{s.t.} \ \hat{\boldsymbol{b}}_1^{\mathrm{H}} \hat{\boldsymbol{b}}_1 = 1. \tag{12}$$

Correspondingly $\boldsymbol{b}_1$ is the eigenvector belonging to the strongest eigenvalue of the channel Gram matrix $\boldsymbol{H}_{\pi(1)} \boldsymbol{H}_{\pi(1)}^{\mathrm{H}}$. As the receive filters are coupled with each other via the pseudoinverse and both the channel gains and the optimum power allocation depend on the pseudoinverse, (11) becomes intractable for $i > 1$. As proposed in [30] for the sum rate maximization with equal weights, we also keep the receive filters of the previous steps fixed, which leads to the following optimization problem:

$$\{\pi(i), \boldsymbol{b}_i\} =$$
$$\underset{\{k, \hat{\boldsymbol{b}}_i\}}{\arg \max} R_{\mathrm{wsr},i}\left(\pi(1), \ldots, \pi(i-1), k, \boldsymbol{b}_1, \ldots, \boldsymbol{b}_{i-1}, \hat{\boldsymbol{b}}_i\right)$$
$$\text{s.t.} \ \hat{\boldsymbol{b}}_i^{\mathrm{H}} \hat{\boldsymbol{b}}_i = 1, \quad \boldsymbol{b}_j^{\mathrm{H}} \hat{\boldsymbol{b}}_i = 0, \ \forall j < i \text{ for } \pi(j) = k. \tag{13}$$

Problem (13) is solved by first computing the optimum receive filters $\boldsymbol{b}_i(k)$ for every user $k = 1, \ldots, K$, i.e.

$$\boldsymbol{b}_i(k) =$$
$$\underset{\hat{\boldsymbol{b}}_i}{\arg \max} R_{\mathrm{wsr},i}\left(\pi(1), \ldots, \pi(i-1), k, \boldsymbol{b}_1, \ldots, \boldsymbol{b}_{i-1}, \hat{\boldsymbol{b}}_i\right), \ \forall k$$
$$\text{s.t.} \ \hat{\boldsymbol{b}}_i^{\mathrm{H}} \hat{\boldsymbol{b}}_i = 1, \boldsymbol{b}_j^{\mathrm{H}} \hat{\boldsymbol{b}}_i = 0, \forall j < i \text{ for } \pi(j) = k. \tag{14}$$

Afterwards these receive filters $\boldsymbol{b}_i(k)$ optimum with respect to (14) are used to compute the weighted sum rates and the next data stream is allocated to that user that leads to the largest weighted sum rate, i.e.,

$$\pi(i) =$$
$$\underset{k}{\arg \max} R_{\mathrm{wsr},i}\left(\pi(1), \ldots, \pi(i-1), k, \boldsymbol{b}_1, \ldots, \boldsymbol{b}_{i-1}, \boldsymbol{b}_i(k)\right). \tag{15}$$

While for given receive filters $\boldsymbol{b}_i(k)$, (15) can be solved straight forwardly via computing (9) with (8) and (10) for every user, (14) is still nonconvex, as all channel gains $\lambda_{i,1}, \ldots, \lambda_{i,i}$ depend on $\hat{\boldsymbol{b}}_i$ via the pseudoinverse [cf. (8)].

*3) Optimization of a Lower Bound:* To make (14) more tractable we use a lower bound for the weighted sum rate in the following. This bound is given by

$$R_{\mathrm{wsr},i}\left(\pi(1),\ldots,\pi(i-1),k,\boldsymbol{b}_1,\ldots,\boldsymbol{b}_{i-1},\hat{\boldsymbol{b}}_i\right)$$

$$\geq \left(\sum_{j=1}^{i-1}\mu_{\pi(j)}+\mu_k\right)\log_2\left(1+\frac{P_{\mathrm{Tx}}}{\left\|\left[\begin{array}{c}\boldsymbol{H}_{\mathrm{co},i-1}\\\hat{\boldsymbol{b}}_i^{\mathrm{H}}\boldsymbol{H}_k\end{array}\right]^+\right\|_{\mathrm{F}}^2}\right) \quad (16)$$

and has been derived in [36]. An alternative proof is presented in Appendix A. Note that (16) is only valid, when all data streams receive powers strictly greater than zero. For the finally chosen user allocation this is always the case, as a data stream with zero power cannot transmit data but imposes zero-forcing constraints on the other users and therefore degrades those users' channel gains. The algorithm has therefore to prevent such a situation by not choosing such a user allocation. As simulation results will reveal, the bound in (16) is tight. Interestingly, this lower bound can be achieved by using the composite channel matrix's pseudoinverse with normalized columns as precoder together with sub-optimum powers $\gamma_{i,j}=P_{\mathrm{Tx}}/(\lambda_{i,j}^2\|\boldsymbol{H}_{\mathrm{co},i}^+\|_{\mathrm{F}}^2)$. This way, each data stream exhibits the same transmission rate. By using the lower bound (16), (14) is approximated by

$$\boldsymbol{b}_i(k)=\arg\min_{\hat{\boldsymbol{b}}_i}\left\|\left[\begin{array}{c}\boldsymbol{H}_{\mathrm{co},i-1}\\\hat{\boldsymbol{b}}_i^{\mathrm{H}}\boldsymbol{H}_k\end{array}\right]^+\right\|_{\mathrm{F}}^2$$

$$\text{s.t. } \hat{\boldsymbol{b}}_i^{\mathrm{H}}\hat{\boldsymbol{b}}_i=1,\quad \boldsymbol{b}_j^{\mathrm{H}}\hat{\boldsymbol{b}}_i=0,\forall j<i \text{ for } \pi(j)=k. \quad (17)$$

Using the successive update of the pseudoinverse with the LQ-decomposition of the matrix $\boldsymbol{H}_{\mathrm{co},i-1}=\boldsymbol{L}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}}$ derived in [27], we have

$$\left\|\left[\begin{array}{c}\boldsymbol{H}_{\mathrm{co},i-1}\\\hat{\boldsymbol{b}}_i^{\mathrm{H}}\boldsymbol{H}_k\end{array}\right]^+\right\|_{\mathrm{F}}^2=\mathrm{tr}\left(\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{L}_{i-1}^{-1}\right)$$

$$+\frac{1+\hat{\boldsymbol{b}}_i^{\mathrm{H}}\boldsymbol{H}_k\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\boldsymbol{H}_k^{\mathrm{H}}\hat{\boldsymbol{b}}_i}{\hat{\boldsymbol{b}}_i^{\mathrm{H}}\boldsymbol{H}_k\left(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\right)\boldsymbol{H}_k^{\mathrm{H}}\hat{\boldsymbol{b}}_i}. \quad (18)$$

As the matrix $\boldsymbol{L}_{i-1}$ is independent of the index $k$ and the filter $\hat{\boldsymbol{b}}_i$, the optimization of the receive filter for user $k$ reduces to [cf. (17)]

$$\boldsymbol{b}_i(k)=\arg\max_{\hat{\boldsymbol{b}}_i}\frac{\hat{\boldsymbol{b}}_i^{\mathrm{H}}\boldsymbol{H}_k\left(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\right)\boldsymbol{H}_k^{\mathrm{H}}\hat{\boldsymbol{b}}_i}{\hat{\boldsymbol{b}}_i^{\mathrm{H}}\left(\boldsymbol{I}_{r_k}+\boldsymbol{H}_k\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\boldsymbol{H}_k^{\mathrm{H}}\right)\hat{\boldsymbol{b}}_i}$$

$$\text{s.t. } \hat{\boldsymbol{b}}_i^{\mathrm{H}}\hat{\boldsymbol{b}}_i=1,\ \boldsymbol{b}_j^{\mathrm{H}}\hat{\boldsymbol{b}}_i=0,\forall j<i \text{ for } \pi(j)=k. \quad (19)$$

where we have inserted the unity norm constraint $1=\hat{\boldsymbol{b}}_i^{\mathrm{H}}\hat{\boldsymbol{b}}_i$ into the denominator of the objective function. The objective

function in (19) is maximized by choosing $\hat{\boldsymbol{b}}_i$ to be a generalized eigenvector belonging to the principal generalized eigenvalue of the matrix pair $\boldsymbol{H}_k(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}})\boldsymbol{H}_k^{\mathrm{H}}$ and $\boldsymbol{I}_{r_k}+\boldsymbol{H}_k\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\boldsymbol{H}_k^{\mathrm{H}}$. As the objective function is independent of the norm of $\hat{\boldsymbol{b}}_i$, the norm-one constraint on $\hat{\boldsymbol{b}}_i$ can be easily fulfilled by taking the norm-one generalized eigenvector. In the following, we will show that additionally the orthogonality constraints $\boldsymbol{b}_j^{\mathrm{H}}\hat{\boldsymbol{b}}_i=0,\forall j<i$ with $\pi(j)=k$ are also fulfilled, when the objective function becomes maximum. At this maximum $\boldsymbol{b}_i(k)$ fulfills

$$\boldsymbol{H}_k\left(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\right)\boldsymbol{H}_k^{\mathrm{H}}\boldsymbol{b}_i(k)$$

$$=\lambda_{k,i-1,\max}\left(\boldsymbol{I}_{r_k}+\boldsymbol{H}_k\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\boldsymbol{H}_k^{\mathrm{H}}\right)\boldsymbol{b}_i(k) \quad (20)$$

where $\lambda_{k,i-1,\max}$ corresponds to the maximum generalized eigenvalue. Note that $\boldsymbol{I}_{r_k}+\boldsymbol{H}_k\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\boldsymbol{H}_k^{\mathrm{H}}$ is always invertible due to the fact that all eigenvalues are greater than or equal to one. Therefore,

$$\boldsymbol{b}_i(k)\in\left(\mathrm{null}\left\{\boldsymbol{H}_k\left(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\right)\boldsymbol{H}_k^{\mathrm{H}}\right\}\right)^{\perp}.$$

Due to the properties of the LQ-decomposition, the vector $\boldsymbol{H}_{\pi(j)}^{\mathrm{H}}\boldsymbol{b}_j=\boldsymbol{H}_k^{\mathrm{H}}\boldsymbol{b}_j=\boldsymbol{Q}_j\boldsymbol{Q}_j^{\mathrm{H}}\boldsymbol{H}_k^{\mathrm{H}}\boldsymbol{b}_j$ lies completely in $\mathrm{range}\{\boldsymbol{Q}_j\boldsymbol{Q}_j^{\mathrm{H}}\}=\mathrm{null}\{(\boldsymbol{I}_N-\boldsymbol{Q}_j\boldsymbol{Q}_j^{\mathrm{H}})\}\subseteq\mathrm{null}\{(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}})\}$. Therefore

$$\boldsymbol{b}_j\in\mathrm{null}\left\{\left(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\right)\boldsymbol{H}_k^{\mathrm{H}}\right\}.$$

As $\boldsymbol{b}_i(k)\in(\mathrm{null}\{\boldsymbol{H}_k(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}})\boldsymbol{H}_k^{\mathrm{H}}\})^{\perp}$, $\boldsymbol{b}_j$ and $\boldsymbol{b}_i(k)$ are orthogonal for all $j<i$ with $\pi(j)=k$. The generalized eigenvalue

$$\lambda_{k,i-1,\max}=\max_{\hat{\boldsymbol{b}}_i}\frac{\hat{\boldsymbol{b}}_i^{\mathrm{H}}\boldsymbol{H}_k\left(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\right)\boldsymbol{H}_k^{\mathrm{H}}\hat{\boldsymbol{b}}_i}{\hat{\boldsymbol{b}}_i^{\mathrm{H}}\left(\boldsymbol{I}_{r_k}+\boldsymbol{H}_k\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\boldsymbol{H}_k^{\mathrm{H}}\right)\hat{\boldsymbol{b}}_i}$$

can be expressed in terms of the maximum eigenvalue of the matrix

$$\boldsymbol{A}_{k,i-1}=\left(\boldsymbol{I}_{r_k}+\boldsymbol{H}_k\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\boldsymbol{H}_k^{\mathrm{H}}\right)^{-1}\boldsymbol{H}_k$$

$$\times\left(\boldsymbol{I}_N-\boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\right)\boldsymbol{H}_k^{\mathrm{H}} \quad (21)$$

i.e.,

$$\lambda_{k,i-1,\max}=\lambda_{\max}(\boldsymbol{A}_{k,i-1}).$$

By using (18) the weighted sum rate $R_{\mathrm{wsr}}(\pi(1),\ldots,\pi(i-1),k,\boldsymbol{b}_1,\ldots,\boldsymbol{b}_{i-1},\boldsymbol{b}_i(k))$ can therefore be also expressed in terms of the maximum eigenvalue of the matrix $\boldsymbol{A}_{k,i-1}$ such that

$$R_{\mathrm{wsr},i}\geq\left(\sum_{j=1}^{i-1}\mu_{\pi(j)}+\mu_k\right)$$

$$\times\log_2\left(1+\frac{P_{\mathrm{Tx}}}{\mathrm{tr}\left(\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{L}_{i-1}^{-1}\right)+\lambda_{\max}^{-1}(\boldsymbol{A}_{k,i-1})}\right). \quad (22)$$

TABLE I
OVERVIEW OF THE DOWNLINK ALGORITHM

**Initialization:**
$$\{\pi(1), \boldsymbol{b}_1\} = \underset{\{k,\hat{\boldsymbol{b}}_1\}}{\arg\max}\, \mu_k \log_2\left(1 + P_{\text{Tx}}\hat{\boldsymbol{b}}_1^{\text{H}}\boldsymbol{H}_k\boldsymbol{H}_k^{\text{H}}\hat{\boldsymbol{b}}_1\right),$$
$$\text{s.t. } \hat{\boldsymbol{b}}_1^{\text{H}}\hat{\boldsymbol{b}}_1 = 1,$$
$$\boldsymbol{H}_{\text{co},1} = \boldsymbol{b}_1^{\text{H}}\boldsymbol{H}_{\pi(1)}, \quad \boldsymbol{B}_{\pi(1)} = \boldsymbol{b}_1, \quad \boldsymbol{B}_k = [\,], \forall k \neq \pi(1),$$
$$R_{\text{wsr,old}} = \log_2\left(1 + P_{\text{Tx}}\boldsymbol{b}_1^{\text{H}}\boldsymbol{H}_{\pi(1)}\boldsymbol{H}_{\pi(1)}^{\text{H}}\boldsymbol{b}_1\right),$$
$$\boldsymbol{Q}_1 = \frac{\boldsymbol{H}_{\pi(1)}^{\text{H}}\boldsymbol{b}_1}{\|\boldsymbol{H}_{\pi(1)}^{\text{H}}\boldsymbol{b}_1\|_2}, \quad \boldsymbol{L}_1^{-1} = \frac{1}{\|\boldsymbol{H}_{\pi(1)}^{\text{H}}\boldsymbol{b}_1\|_2}, \quad i = 2$$
**while** $i \leq N$ **do**
$$\pi(i) = \underset{k}{\arg\max}\, R_{\text{wsr},i}\left(\pi(1),\ldots,\pi(i-1),k,\boldsymbol{b}_1,\ldots,\boldsymbol{b}_{i-1},\boldsymbol{b}_i(k)\right)$$
$$\boldsymbol{b}_i(k) = \underset{\hat{\boldsymbol{b}}_i}{\arg\max}\, f_{\text{rx}}\left(\hat{\boldsymbol{b}}_i,\boldsymbol{H}_k,\boldsymbol{H}_{\text{co},i-1}\right), \quad \text{s.t. } \hat{\boldsymbol{b}}_i^{\text{H}}\hat{\boldsymbol{b}}_i = 1,$$
$$f_{\text{rx}}(\hat{\boldsymbol{b}}_i,\boldsymbol{H}_k,\boldsymbol{H}_{\text{co},i-1}) = \frac{\hat{\boldsymbol{b}}_i^{\text{H}}\boldsymbol{H}_k(\boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\text{H}})\boldsymbol{H}_k^{\text{H}}\hat{\boldsymbol{b}}_i}{\hat{\boldsymbol{b}}_i^{\text{H}}\left(\boldsymbol{I}_{r_k} + \boldsymbol{H}_k\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\text{H}}\boldsymbol{Q}_{i-1}^{\text{H}}\boldsymbol{H}_k^{\text{H}}\right)\hat{\boldsymbol{b}}_i},$$
$$\text{or } f_{\text{rxfilter}}\left(\hat{\boldsymbol{b}}_i,\boldsymbol{H}_k,\boldsymbol{H}_{\text{co},i-1}\right) = \hat{\boldsymbol{b}}_i^{\text{H}}\boldsymbol{H}_k(\boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\text{H}})\boldsymbol{H}_k^{\text{H}}\hat{\boldsymbol{b}}_i$$
$$R_{\text{wsr,new}} = R_{\text{wsr},i}\left(\pi(1),\ldots,\pi(i),\boldsymbol{b}_1,\ldots,\boldsymbol{b}_{i-1},\boldsymbol{b}_i(\pi(i))\right)$$
**if** $R_{\text{wsr,new}} > R_{\text{wsr,old}}$ **then**
$$\boldsymbol{b}_i = \boldsymbol{b}_i(\pi_i), \quad \boldsymbol{B}_{\pi(i)} = \left[\boldsymbol{B}_{\pi(i)} \,\, \boldsymbol{b}_i\right]$$
$$\boldsymbol{H}_{\text{co},i} = \begin{bmatrix} \boldsymbol{H}_{\text{co},i-1} \\ \boldsymbol{b}_i^{\text{H}}\boldsymbol{H}_{\pi(i)} \end{bmatrix},$$
$$\boldsymbol{Q}_i = \left[\boldsymbol{Q}_{i-1} \,\, \frac{(\boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\text{H}})\boldsymbol{H}_{\pi(i)}^{\text{H}}\boldsymbol{b}_i}{\ell_{i,i}}\right],$$
$$\boldsymbol{L}_i^{-1} = \begin{bmatrix} \boldsymbol{L}_{i-1}^{-1} & \boldsymbol{0}_{i-1,1} \\ -\frac{\boldsymbol{b}_i^{\text{H}}\boldsymbol{H}_{\pi(i)}\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}}{\left\|(\boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\text{H}})\boldsymbol{H}_{\pi(i)}^{\text{H}}\boldsymbol{b}_i\right\|_2} & \frac{1}{\ell_{i,i}} \end{bmatrix},$$
$$\ell_{i,i} = \left\|\left(\boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\text{H}}\right)\boldsymbol{H}_{\pi(i)}^{\text{H}}\boldsymbol{b}_i\right\|_2$$
$$R_{\text{wsr,old}} = R_{\text{wsr,new}}, \quad i = i + 1$$
**else**
$$i = i - 1, \text{ break}$$
**end if**
**end while**
$\boldsymbol{\Psi}_{\text{eff},i} = [\boldsymbol{\psi}_1 \ldots \boldsymbol{\psi}_i] = \boldsymbol{H}_{\text{co},i}^+ \boldsymbol{\Lambda}_i$, $\boldsymbol{\Lambda}_i = \text{diag}(\lambda_{i,1},\ldots,\lambda_{i,i})$, scales
columns of $\boldsymbol{H}_{\text{co},i}^+$ to unit norm
$\gamma_{i,j} \leftarrow$ waterfilling of $P_{\text{Tx}}$ over $\lambda_{i,1},\ldots,\lambda_{i,i}$
$\boldsymbol{P}_k = [\,], \quad \forall k$
**for** $j = 1$ to $i$ **do**
$$\boldsymbol{P}_{\pi(j)} = \left[\boldsymbol{P}_{\pi(j)} \,\, \boldsymbol{\psi}_j\sqrt{\gamma_{i,j}}\right]$$
**end for**

In summary we have made the following simplifications to approximate the optimum solution to (2) at drastically reduced computational complexity. At first, zero-forcing constraints have been introduced between data streams of different users as well as data streams allocated to the same user. Secondly, the allocation of data streams to users and the determination of the corresponding receive filters has been conducted in a successive manner, where in each step only the next allocated user and its receive filter are determined. Thirdly, the receive filters have been chosen to maximize a lower bound for the weighted sum rate.

## IV. DOWNLINK ALGORITHM

In this section we will give an overview of the proposed algorithm operating directly in the downlink and comment on some implementation aspects. A pseudocode of the algorithm is given in Table I.

### A. Initialization

The first user and the corresponding filter are determined using (12), i.e.,

$$\{\pi(1), \boldsymbol{b}_1\} = \underset{\{k,\hat{\boldsymbol{b}}_1\}}{\arg\max}\, \mu_k \log_2\left(1 + P_{\text{Tx}}\hat{\boldsymbol{b}}_1^{\text{H}}\boldsymbol{H}_k\boldsymbol{H}_k^{\text{H}}\hat{\boldsymbol{b}}_1\right)$$
$$\text{s.t. } \hat{\boldsymbol{b}}_1^{\text{H}}\hat{\boldsymbol{b}}_1 = 1.$$

### B. Successive Allocation of Data Streams

The proposed method assumes that the users allocated in the previous steps and the corresponding receive filters are kept fixed and optimize the newly allocated data stream and the corresponding receive filter at the corresponding terminal, as described in the following for a certain step $i$.

*1) Determination of Filters at the Terminals:* For each user $k$, whose allocated number of data streams is less than its number of antennas, the receive filter $\boldsymbol{b}_i(k)$ in the downlink maximizing a lower bound for the weighted sum rate is computed according to (19)

$$\boldsymbol{b}_i(k) = \underset{\hat{\boldsymbol{b}}_i}{\arg\max}\, \frac{\hat{\boldsymbol{b}}_i^{\text{H}}\boldsymbol{H}_k\left(\boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\text{H}}\right)\boldsymbol{H}_k^{\text{H}}\hat{\boldsymbol{b}}_i}{\hat{\boldsymbol{b}}_i^{\text{H}}\left(\boldsymbol{I}_{r_k} + \boldsymbol{H}_k\boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\text{H}}\boldsymbol{Q}_{i-1}^{\text{H}}\boldsymbol{H}_k^{\text{H}}\right)\hat{\boldsymbol{b}}_i}$$
$$\text{s.t. } \hat{\boldsymbol{b}}_i^{\text{H}}\hat{\boldsymbol{b}}_i = 1 \tag{23}$$

where we have omitted the orthogonality constraints of (19), which are inactive at the optimum. Note that above optimization is independent of the weights $\mu_1,\ldots,\mu_K$. The same equation has already been derived for the special case of equal weights for all users in the downlink in [30] and [31]. Hence, we have shown that the same receive filters also maximize a lower bound for the weighted sum rate. For the sake of complexity reduction the matrix inversion necessary in (23) can be avoided at the expense of slight performance losses by using different downlink filters

$$\tilde{\boldsymbol{b}}_i(k) = \underset{\hat{\boldsymbol{b}}_i}{\arg\max}\, \hat{\boldsymbol{b}}_i^{\text{H}}\boldsymbol{H}_k\left(\boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\text{H}}\right)\boldsymbol{H}_k^{\text{H}}\hat{\boldsymbol{b}}_i$$
$$\text{s.t. } \hat{\boldsymbol{b}}_i^{\text{H}}\hat{\boldsymbol{b}}_i = 1. \tag{24}$$

Using the same reasoning as in [30] and [31] it can be shown that the $\tilde{\boldsymbol{b}}_i(k)$ maximize an even looser lower bound for the weighted sum rate. A derivation of this bound is given in Appendix B.

*2) User Selection:* For the user selection, the next data stream is tentatively allocated to each user applying the receive filters $\boldsymbol{b}_i(k)$. For each allocation, the resulting weighted sum rate is computed and the user that leads to the strongest increase in weighted sum rate is finally selected, i.e.,

$$\pi(i) = \underset{k}{\arg\max}\, R_{\text{wsr},i}\left(\pi(1),\ldots,\pi(i-1),k,\right.$$
$$\left.\boldsymbol{b}_1,\ldots,\boldsymbol{b}_{i-1},\boldsymbol{b}_i(k)\right). \tag{25}$$

In case of the simplified receive filter determination (24), the corresponding receive filters $\tilde{\boldsymbol{b}}_1,\ldots,\tilde{\boldsymbol{b}}_{i-1},\tilde{\boldsymbol{b}}_i(k)$ need to be used in (25). Consequently a different user allocation can be obtained. If it is observed that zero power is allocated to a data stream, this data stream must be removed. This removal will lead most likely to an increase of the other subchannel gains and therefore to an

increase in weighted sum rate as one zero-forcing constraint is removed from the optimization problem.

The user selection can be alternatively performed with the lower bound from (22), i.e.

$$\pi(i) = \arg\max_k \left( \sum_{j=1}^{i-1} \mu_{\pi(j)} + \mu_k \right) \times \log_2 \left( 1 + \frac{P_{\mathrm{Tx}}}{\mathrm{tr}\left( \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{L}_{i-1}^{-1} \right) + \lambda_{\max}^{-1}(\boldsymbol{A}_{k,i-1})} \right).$$

The usage of this bound only leads to marginal gains in terms of computational complexity, since the power allocation according to (10) exhibits negligible computational complexity compared to the computation of the principal generalized eigenvalues. However, this bound can be used for a user preselection. This preselection implies that some users are excluded without explicitly computing their generalized eigenvalues, which helps to reduce the computational complexity dramatically. As shown in Appendix C, the maximum eigenvalue of the matrix $\boldsymbol{A}_{k,i}$ can be bounded as follows:

$$\lambda_{\mathrm{lb},k,i} = \frac{\mathrm{tr}\left( \boldsymbol{H}_k \left( \boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}} \right) \boldsymbol{H}_k^{\mathrm{H}} \right)}{r_k \left( 1 + \mathrm{tr}\left( \boldsymbol{H}_k \boldsymbol{Q}_{i-1}\boldsymbol{L}_{i-1}^{-1}\boldsymbol{L}_{i-1}^{-\mathrm{H}}\boldsymbol{Q}_{i-1}^{\mathrm{H}}\boldsymbol{H}_k^{\mathrm{H}} \right) \right)}$$

$$\leq \lambda_{\max}(\boldsymbol{A}_{k,i-1}) \leq \mathrm{tr}\left( \boldsymbol{H}_k \left( \boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}} \right) \boldsymbol{H}_k^{\mathrm{H}} \right)$$

$$= \lambda_{\mathrm{ub},k,i}. \tag{26}$$

The rationale of the proposed user preselection grounds on the conservative rule that, if the upper bound of the estimated sum rate allocating the next data stream to a certain user is smaller than the maximum lower bound over all user allocations, that user will certainly not lead to the maximum increase in estimated weighted sum rate. Hence, only for users $m$ with

$$\left( \sum_{j=1}^{i-1} \mu_{\pi(j)} + \mu_m \right) \log_2 \left( 1 + \frac{P_{\mathrm{Tx}}}{\mathrm{tr}\left( \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{L}_{i-1}^{-1} \right) + \frac{1}{\lambda_{\mathrm{ub},m,i}}} \right)$$

$$> \max_k \left( \sum_{j=1}^{i-1} \mu_{\pi(j)} + \mu_k \right)$$

$$\times \log_2 \left( 1 + \frac{P_{\mathrm{Tx}}}{\mathrm{tr}\left( \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{L}_{i-1}^{-1} \right) + \frac{1}{\lambda_{\mathrm{lb},k,i}}} \right) \tag{27}$$

the maximum generalized eigenvalue needs to be computed. Thus, with almost no additional effort matrix inversions can be avoided and therefore the total complexity of the algorithm is reduced drastically. Although this rule might be considered as too conservative it turns out that in practical applications it has a considerable impact on complexity reduction without any loss of performance. For the selection of the first data stream according to (12), which requires the computation of the principal eigenvalues of the matrices $\boldsymbol{H}_k\boldsymbol{H}_k^{\mathrm{H}}$, the lower and upper bounds are given by

$$\lambda_{\mathrm{lb},k,1} = \frac{\mathrm{tr}\left( \boldsymbol{H}_k\boldsymbol{H}_k^{\mathrm{H}} \right)}{r_k}, \quad \lambda_{\mathrm{ub},k,1} = \mathrm{tr}\left( \boldsymbol{H}_k\boldsymbol{H}_k^{\mathrm{H}} \right)$$
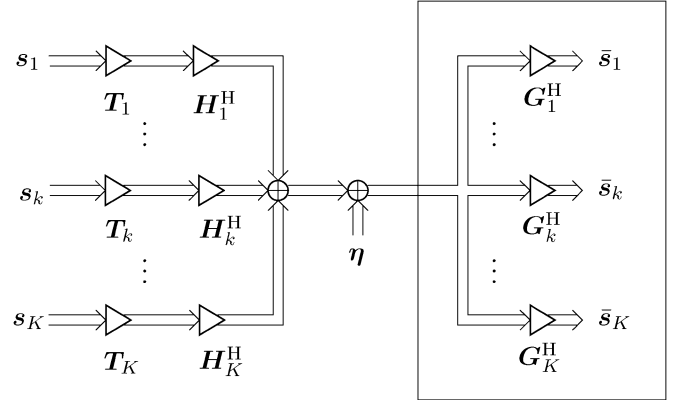


Fig. 2. System model of the MIMO multiple access channel.

which coincides with the bounds for principal eigenvalues in [37, Ch 2.3] and leads to the same user preselection as proposed in [38].

### C. Termination

As with each next allocated data stream, the channel gains of the previously allocated subchannels decrease, it might happen that all possible new data stream allocations do not lead to any increase in weighted sum rate. In that case the algorithm must be terminated, which happens at the latest after $i = N$.

### D. Signaling of Receive Filters

Determining the filters $\boldsymbol{b}_i$ at the terminals directly would require the knowledge of all users' channel matrices at all terminals, which is rather unrealistic in practice. The $\boldsymbol{b}_i$ are therefore determined solely at the transmitter. In a signaling phase before data transmission these filters are made known to the receivers together with the user allocation. Similarly to [32] or [39], this can be done as follows. First common pilot symbols are sent to the users, where the pilot symbols are precoded such that the estimate is equal to the filters to be applied at the corresponding users. In a second step user identifiers are sent over the resulting subchannels such that each user is able to detect on which subchannels he will receive data. Alternatively, [40] proposes signaling schemes with quantized feedforward of the terminals' filters. It should be noted that the problem of signaling the receive filters arises generally in the MIMO broadcast channel with perfect CSI at the transmitter, i.e., also with the algorithms from [15], [16], [20], or [21].

## V. DUALITY BASED WEIGHTED SUM RATE MAXIMIZATION

In this section we will present an algorithm that performs the successive data stream allocation and filter determination in the dual uplink and leads to additional gains by using minimum mean square error (MMSE) filters in the uplink and further optimizing the receive filters after the transformation in the downlink. First the system model of the dual uplink is introduced, which is depicted in Fig. 2.

In the dual multiple access channel, $K$ decentralized users transmit their data to the base station. To this end, user $k$'s symbol vector $\boldsymbol{s}_k$ is precoded with the matrix $\boldsymbol{T}_k \in \mathbb{C}^{r_k \times d_k}$ and the filtered symbol vector then propagates over the Hermitian

channel $\boldsymbol{H}_k^{\mathrm{H}}$. At the receiver side, the signals of the $K$ users are summed up and zero mean circularly symmetric Gaussian noise $\boldsymbol{\eta} \in \mathbb{C}^N$ with identity covariance matrix is added. Finally, the linear filter $\boldsymbol{G}_k^{\mathrm{H}} \in \mathbb{C}^{d_k \times N}$ generates the respective symbol estimate $\bar{\boldsymbol{s}}_k \in \mathbb{C}^{d_k}$ of user $k$ out of the noisy received signal, with $k$ ranging from 1 to $K$. Therefore, the symbol estimate $\bar{\boldsymbol{s}}_k$ in the dual uplink channel reads as

$$\bar{\boldsymbol{s}}_k = \boldsymbol{G}_k^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \boldsymbol{T}_k \boldsymbol{s}_k + \boldsymbol{G}_k^{\mathrm{H}} \sum_{\substack{j=1 \\ j \neq k}}^{K} \boldsymbol{H}_j^{\mathrm{H}} \boldsymbol{T}_j \boldsymbol{s}_j + \boldsymbol{G}_k^{\mathrm{H}} \boldsymbol{\eta}.$$

Similarly to the downlink rate term in (1), the rate of user $k$ that can be achieved in the dual uplink under linear filtering can be expressed as

$$R_{k,\mathrm{UL}} = \log_2 \frac{\left| \boldsymbol{G}_k^{\mathrm{H}} \boldsymbol{G}_k + \sum_{j=1}^{K} \boldsymbol{G}_k^{\mathrm{H}} \boldsymbol{H}_j^{\mathrm{H}} \boldsymbol{T}_j \boldsymbol{T}_j^{\mathrm{H}} \boldsymbol{H}_j \boldsymbol{G}_k \right|}{\left| \boldsymbol{G}_k^{\mathrm{H}} \boldsymbol{G}_k + \sum_{\substack{j=1 \\ j \neq k}}^{K} \boldsymbol{G}_k^{\mathrm{H}} \boldsymbol{H}_j^{\mathrm{H}} \boldsymbol{T}_j \boldsymbol{T}_j^{\mathrm{H}} \boldsymbol{H}_j \boldsymbol{G}_k \right|}.$$

In the dual uplink, we can optimize the weighted sum rate under a total power constraint $\sum_{k=1}^{K} \mathrm{tr}(\boldsymbol{T}_k \boldsymbol{T}_k^{\mathrm{H}}) \leq P_{\mathrm{Tx}}$ in the same successive fashion as in the original downlink. This results from the rate region duality between the MAC and the BC under linear filtering, since any rate tuple in the broadcast channel can also be achieved in the dual multiple access channel, and vice versa, see [19]. Following the same argumentation as in the downlink, jointly optimized transmit and receive filters cannot easily be obtained one by one. Therefore, we restrict the filters to completely suppress the interstream interference between all allocated streams and end up with a zero-forcing system in the dual MAC as well. Additionally, we again operate on a lower bound of the weighted sum rate instead of the original utility.

Let $\boldsymbol{u}_i$ denote the unit-norm beamforming vector associated to the data stream that is allocated at the $i$th step of the successive algorithm and that belongs to user $\pi(i)$. To obtain $\boldsymbol{t}_i$, the vector $\boldsymbol{u}_i$ must be multiplied by the square root of the power allocated to the $i$th data stream[2]. Then, we can define a composite channel matrix for the $i$th stage in the dual MAC as the concatenation of the channel and the unit-norm beamformers

$$\boldsymbol{H}_{\mathrm{co},i}^{\mathrm{H}} = \left[ \boldsymbol{H}_{\pi(1)}^{\mathrm{H}} \boldsymbol{u}_1, \ldots, \boldsymbol{H}_{\pi(i)}^{\mathrm{H}} \boldsymbol{u}_i \right] \in \mathbb{C}^{N \times i}. \tag{28}$$

Storing the receive filters $\boldsymbol{g}_1^{\mathrm{H}}, \ldots, \boldsymbol{g}_i^{\mathrm{H}}$ for the $i$ active data streams in the effective receive filter of stage $i$

$$\boldsymbol{G}_{\mathrm{eff},i}^{\mathrm{H}} = \begin{bmatrix} \boldsymbol{g}_1^{\mathrm{H}} \\ \vdots \\ \boldsymbol{g}_i^{\mathrm{H}} \end{bmatrix} \in \mathbb{C}^{i \times N} \tag{29}$$

the zero-forcing condition can compactly be written as

$$\boldsymbol{G}_{\mathrm{eff},i}^{\mathrm{H}} \boldsymbol{H}_{\mathrm{co},i}^{\mathrm{H}} = \mathbf{I}_i \tag{30}$$

and leads to the effective receive filters

$$\boldsymbol{G}_{\mathrm{eff},i}^{\mathrm{H}} = \boldsymbol{H}_{\mathrm{co},i}^{\mathrm{H},+} \tag{31}$$

representing the pseudoinverse of the composite channel matrix. Feeding the unit-norm beamformers $\boldsymbol{u}_1, \ldots, \boldsymbol{u}_i$ with powers $\gamma_{i,1}, \ldots, \gamma_{i,i}$ at stage $i$, the received signal-to-noise ratio (SNR) of the $j$th stream with $j \leq i$ reads as

$$\mathrm{SNR}_{i,j} = \frac{\gamma_{i,j}}{\|\boldsymbol{g}_j\|_2^2} = \frac{\gamma_{i,j}}{\left\| \boldsymbol{H}_{\mathrm{co},i}^+ \boldsymbol{e}_j \right\|_2^2} = \gamma_{i,j} \lambda_{i,j}^2 \tag{32}$$

cf. (8) for the downlink counterpart. Using (32), the weighted sum-rate which is obtained by the zero-forcing transmission strategy can be expressed as

$$R_{\mathrm{wsr},i}\left(\pi(1), \ldots, \pi(i), \boldsymbol{u}_1, \ldots, \boldsymbol{u}_i\right)$$
$$= \sum_{j=1}^{i} \mu_{\pi(j)} \log_2\left(1 + \gamma_{i,j} \lambda_{i,j}^2\right) \tag{33}$$

and completely coincides with the downlink term (9) up to the different variable names of the arguments. Note that the optimum water-filling results from (10) in the uplink as well. The key observation here is that from now on, the same successive algorithm to determine the receive filters $\boldsymbol{b}_1, \ldots, \boldsymbol{b}_i$ in the downlink can be applied to determine the unit-norm beamformers $\boldsymbol{u}_1, \ldots, \boldsymbol{u}_i$ in the dual uplink, no matter how this algorithm selects those filters in detail. The power allocation does not differ in both domains and the same SNRs are obtained for the active data streams. The only difference is that any successive algorithm in the downlink returns the *receive* filters from which the transmitters are found via the pseudoinverse of the composite channel, while the same successive algorithm in the dual uplink returns the *transmit* filters.

It is worth to mention that due to the fact that the dual MAC achieves the same rates as the original broadcast channel if transmitters and receivers of both domains are designed according to the successive algorithm, we would not have gained anything by designing the successive algorithm in the dual uplink. In order to go beyond this result, we propose replacing the receive filters that arise from the zero-forcing algorithm through optimum MMSE receive filters such that the joint decoding of all streams of every single user leads to an increase in every user's rate. Doing so, the dual MAC system reacts differently on the application of MMSE receivers

$$\boldsymbol{G}_k^{\mathrm{H}} = \boldsymbol{T}_k^{\mathrm{H}} \boldsymbol{H}_k \left( \mathbf{I}_N + \sum_{\ell=1}^{K} \boldsymbol{H}_\ell^{\mathrm{H}} \boldsymbol{T}_\ell \boldsymbol{T}_\ell^{\mathrm{H}} \boldsymbol{H}_\ell \right)^{-1} \tag{34}$$

than the original system in the broadcast channel on its MMSE receivers

$$\boldsymbol{B}_k^{\mathrm{H}} = \boldsymbol{P}_k^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \left( \mathbf{I}_{r_{qk}} + \boldsymbol{H}_k \sum_{\ell=1}^{K} \boldsymbol{P}_\ell \boldsymbol{P}_\ell^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \right)^{-1}. \tag{35}$$

---

[2]The filters $\boldsymbol{T}_k$ are obtained from the $\boldsymbol{t}_i$ in the same way as the filters $\boldsymbol{B}_k$ are determined from the $\boldsymbol{b}_i$ in the downlink.
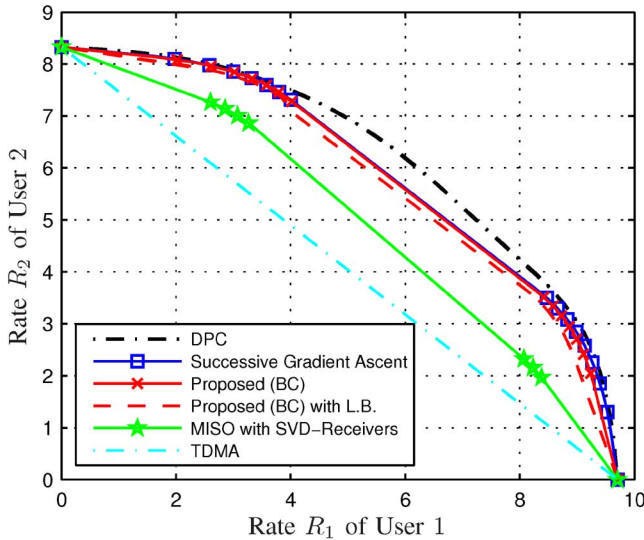
Fig. 3. Supported rate regions for different algorithms. $K = 2$ three-antenna users are served by an $N = 4$ antenna base station. The weights $\mu_1$ and $\mu_2$ have been varied to achieve the complete region.

On average, MMSE receivers lead to a larger increase of the rates in the dual MAC, but not for every single channel realization. Moreover, the transformation from the dual MAC to the downlink entails the possibility for an additional gain of rates. Since the duality transformation generates receive filters in the broadcast channel which preserve the obtained rates in the dual MAC but are not optimum, we can furthermore replace the receivers in the transformed broadcast channel by MMSE receivers. Summing up, we first perform the successive data stream allocation with zero-forcing constraints until no increase in weighted sum rate can be achieved by a further allocation. Then we equip the dual MAC with MMSE receivers instead of the ones of the algorithm in (31). This could also be done directly in the broadcast channel but leads to smaller gains there on average. Next, we convert the obtained filter pairs $\boldsymbol{T}_1, \boldsymbol{G}_1^{\mathrm{H}}, \ldots, \boldsymbol{T}_K, \boldsymbol{G}_K^{\mathrm{H}}$ to the downlink. The arising precoders $\boldsymbol{P}_1, \ldots, \boldsymbol{P}_k$ are then used to set up the MMSE receivers via (35). Due to the special properties of our rate duality, the joint decoding of every user's streams will not become necessary, separate stream decoding also achieves the two mentioned rate gains, see [19]. As in the previous section, all filters are computed at the base station and then signaled to the corresponding user terminals.

## VI. SIMULATION RESULTS

In order to evaluate the performance of the proposed algorithms, we first plot the achievable two user rate regions by varying the user weights $\mu_1$ and $\mu_2$ for a particular channel realization and a fixed transmit power. Afterwards, we present ergodic results by averaging the weighted sum rates over many channel realizations for fixed weights and different transmit powers. In Fig. 3, the rate $R_2$ of user 2 is plotted versus the rate $R_1$ of user 1 where both users have $r_k = 3$ antennas each and are served by an $N = 4$ antenna base station. The chosen channel realization has been drawn from a complex Gaussian

distribution with independently and identically distributed entries, i.e., $\mathrm{vec}([\boldsymbol{H}_1^{\mathrm{T}}, \ldots, \boldsymbol{H}_K^{\mathrm{T}}]) \sim \mathcal{CN}(\boldsymbol{0}, \boldsymbol{I})$. Obviously, all rate regions of the linear filters are outer bounded by the capacity region of the precoder incorporating nonlinear interference cancellation via dirty paper coding. This outer bound (*DPC*) corresponds to the black dash-dotted curve. Rate pairs on the boundary of this region have been obtained with the weighted sum rate maximization algorithm from [3] and the connection of widely separated rate pairs is done via time sharing. Switching to the class of linear precoders, the successive stream allocation approach in [15] utilizing iterative projected gradient ascent steps performs best. Plotted with blue square markers (*Successive Gradient Ascent*), its performance comes at the expense of a relatively high computational complexity compared to the two schemes proposed in this paper. These two have approximately the same performance as the iterative algorithm in [15], irrespective of whether they are implemented in the broadcast channel directly (*Proposed BC*, red curve, cross marker) or in the dual MAC with the twofold MMSE gain, but much smaller complexity. As the rate region of the algorithm operating directly in the broadcast exhibits hardly any visible difference with the MAC+duality algorithm in this scenario, we have omitted the latter region for a better visualization. If a looser lower bound on the weighted sum rate is maximized by choosing the precoders according to (24) to avoid the matrix inversion, the dashed red curve (*Proposed (BC) with L.B.*) is obtained which has to face slight deteriorations. The multiuser MIMO system can be converted to a multiuser MISO system model with virtual users by choosing the left singular vectors of the channel matrices as the receive filters for the respective users [28], [29]. Doing so, the green curve with the star markers (*MISO with SVD-Receivers*) is achieved. This algorithm is only slightly less complex than the proposed method. Although it requires no generalized eigenvalue computations, a complete SVD has to be computed for each user and instead of $K$ users in each step $r = \sum_{k=1}^{K} r_k$ virtual users need to be tested for the maximum weighted sum rate, which implies that for each of those virtual users the precoders have to be computed. Furthermore it offers no possibilities for complexity reduction, as the proposed method. By using the further lower bound for the determination of filters at the terminals a better performance than with the MISO algorithm can be achieved at a lower computational complexity. Finally, the cyan-colored dash-dotted curve (*TDMA*) shows the rate region of the time division multiple access scheme. We did not simulate the algorithm in [16] for weighted sum rate maximization due to complexity reasons. In its original form in [16], the objective can only be converted to a posynomial if all users have rational weights. In this case, expanding the objective would lead to an huge number of monomials in the posynomial. Moreover, the geometric programming toolbox will rather apply the log operator to the objective instead of expanding it. Thus, [16] inherently corresponds to solving the original weighted sum rate problem where every summand corresponds to the weighted rate of the respective user.

To overcome the influence of a particular channel realization, we average over 1000 independent realizations and increase the number of users from two to $K = 4$, again with
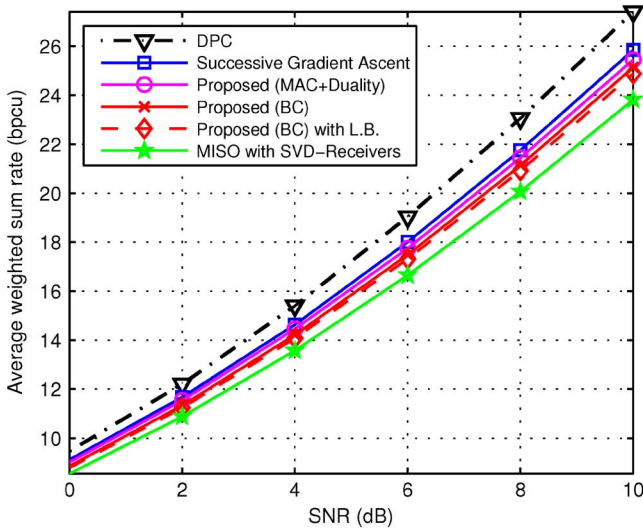
Fig. 4. Weighted sum rate versus transmit power for six different algorithms. $K = 4$ three-antenna users are served by an $N = 4$ antenna base station. The weights are $\mu_1 = \mu_2 = 2$ and $\mu_3 = \mu_4 = 1$.
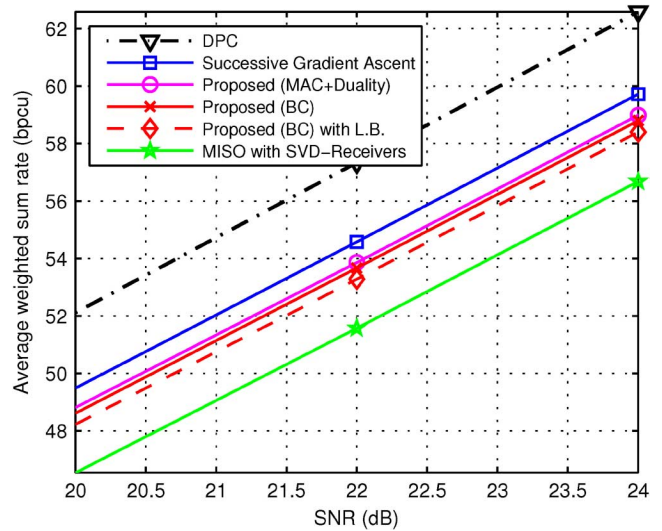


Fig. 5. Weighted sum rate versus transmit power for six different algorithms. $K = 4$ three-antenna users are served by an $N = 4$ antenna base station. The weights are $\mu_1 = \mu_2 = 2$ and $\mu_3 = \mu_4 = 1$.

$r_k = 3$ antennas each. Moreover, the user weights read as $\mu_1 = \mu_2 = 2$ and $\mu_3 = \mu_4 = 1$. Fig. 4 shows the averaged weighted sum rate versus the logarithmic transmit power in the small power regime. For very little transmit power $P_{\mathrm{Tx}}$, all curves will coincide as only one data stream will be allocated according to the best metric $\mu_k \log_2(1 + P_{\mathrm{Tx}}\lambda_{\max}(\boldsymbol{H}_k\boldsymbol{H}_k^{\mathrm{H}}))$, see Section IV-A. If the transmit power is increased, different weighted sum rates are obtained by the simulated algorithms since now, more than one data stream will be activated. The performance ranking of all algorithms in Fig. 4 matches that from Fig. 3, and we again observe, that the two proposed complexity reduced algorithms (*Proposed (MAC+Duality)* and *Proposed (BC)*) almost achieve the same weighted sum rate as the iterative algorithm (*Successive Gradient Ascent*) from [15] despite their noniterative structure. Using the looser bound on the weighted sum rate according to (24) to avoid a matrix inversion for the computation of the strongest eigenmode (*Proposed (BC) with L.B.*) almost does not change the performance, which is still better than the one of the MISO algorithm with right singular vectors of the channels as receive filters (*MISO with SVD-Receivers*), see [28] and [29]. Simulation results for transmit powers between 20 and 24 dB are shown in Fig. 5. All curves now seem to have the same slope and differ only by a power shift. Again, the performance ranking does not change even for such high $P_{\mathrm{Tx}}$. For the chosen user weights, the two proposed algorithms reach the weighted sum rate of the best simulated linear algorithm up to a power shift of only 0.15 and 0.2 dB, respectively. By applying the user preselection from (27), 25% of the necessary generalized eigenvalue computations and 50% of the eigenvalue computations necessary for the selection of the user to which the first data stream is allocated to can be avoided in this setting at 24 dB at no performance losses. This reduced complexity together with the little performance losses clearly motivates the application of the two proposed algorithms for weighted sum rate maximization under linear filtering. For the channel averaged
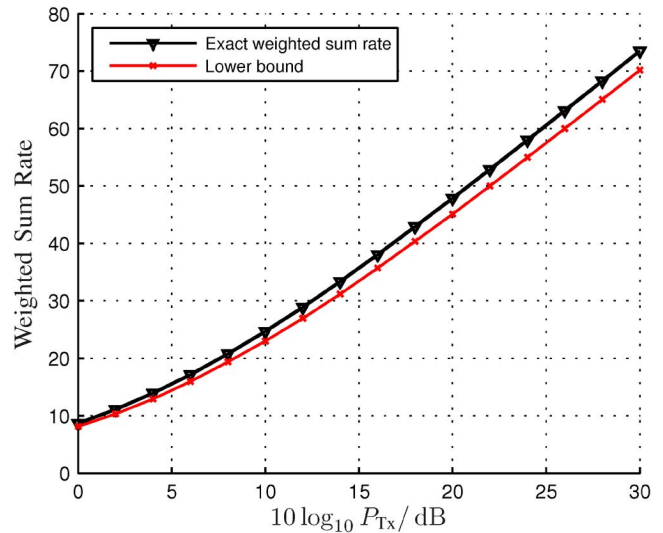


Fig. 6. Actual and estimated weighted sum rate versus transmit power for the proposed algorithm in the broadcast channel. $K = 4$ three-antenna users are served by an $N = 4$ antenna base station. The weights are $\mu_1 = \mu_2 = 2$ and $\mu_3 = \mu_4 = 1$.

sum rate results the simulation of the algorithm in [16] has again prohibitive complexity. Given the weights $\mu_1 = \mu_2 = 2$ and $\mu_3 = \mu_4 = 1$, the objective in [16, (9)] has more than $10^{20}$ summands.

Finally, Fig. 6 compares the average weighted sum rate achievable with the proposed algorithm in the broadcast with the lower bound derived in (16). The same setting as for the previous two simulations has been chosen, i.e., there are $K = 4$ users each equipped with $r_k = 3$ antennas and a base station with $N = 4$ antennas. The weights have been chosen according to $\mu_1 = \mu_2 = 2$, $\mu_3 = \mu_4 = 1$. The weighted sum rate and the estimated sum rate has been averaged over 1000 independent channel realizations. As shown in Fig. 6, the derived bound approaches the actual weighted sum rate quite well.

## VII. Conclusion

In this paper, we have presented two heuristic algorithms which aim at maximizing a weighted sum rate in the MIMO broadcast channel by means of linear signal processing. The first algorithm directly operates in the downlink, while the other one works in the dual uplink, from which the downlink solutions can be obtained via a general rate duality. To avoid the non-convex and combinatorial optimization problems zero-forcing constraints are introduced and a successive approach is used in both cases. In each step, the user to which the next data stream is allocated to and its corresponding filter are determined, such that the increase in weighted sum rate becomes maximum, whereas the user allocation and filters from previous steps are kept fixed. To simplify things further a lower bound to compute the weighted sum rate has been applied. This way, the filters at the terminals can be determined by a generalized eigenvalue problem.

## Appendix A
### Derivation of a Lower Bound for the Weighted Sum Rate

Plugging the optimum power allocation from (10) into the weighted sum rate (9) and exchanging the sum of logarithms by the logarithm of the product leads to

$$R_{\mathrm{wsr},i} = \log_2 \prod_{j=1}^{i} \left( \frac{\mu_{\pi(j)}}{\sum\limits_{\ell=1}^{i} \mu_{\pi(\ell)}} \left( P_{\mathrm{Tx}} + \sum_{\ell=1}^{i} \frac{1}{\lambda_{i,\ell}^2} \right) \lambda_{i,j}^2 \right)^{\mu_{\pi(j)}}$$

where the Lagrange multiplier in (10) reads as

$$\eta_i = \frac{1}{\sum_{\ell=1}^{i} \mu_{\pi(\ell)}} \left( P_{\mathrm{Tx}} + \sum_{\ell=1}^{i} \frac{1}{\lambda_{i,\ell}^2} \right)$$

under the assumption that all powers $\gamma_{i,j}$ are strictly greater than zero. As already explained below (16), for the finally chosen user allocation this assumption always holds. By applying the inequality between the weighted geometric and the weighted

harmonic mean (e.g., [33, Lemma 1]), we obtain a lower bound for the weighted sum rate:

$$
\begin{aligned}
R_{\mathrm{wsr},i} \\
&= \left( \sum_{\ell=1}^{i} \mu_{\pi(\ell)} \right) \\
&\quad \times \log_2 \prod_{j=1}^{i} \left( \frac{\mu_{\pi(j)}}{\sum_{\ell=1}^{i} \mu_{\pi(\ell)}} \right. \\
&\qquad\qquad \left. \times \left( P_{\mathrm{Tx}} + \sum_{\ell=1}^{i} \frac{1}{\lambda_{i,\ell}^2} \right) \lambda_{i,j}^2 \right)^{\frac{\mu_{\pi(j)}}{\sum_{\ell=1}^{i} \mu_{\pi(\ell)}}} \\
&\geq \left( \sum_{\ell=1}^{i} \mu_{\pi(\ell)} \right) \\
&\quad \times \log_2 \frac{\sum_{\ell=1}^{i} \mu_{\pi(\ell)}}{\sum_{j=1}^{i} \mu_{\pi(j)} \left( \frac{\mu_{\pi(j)} \lambda_{i,j}^2}{\sum_{\ell=1}^{i} \mu_{\pi(\ell)}} \left( P_{\mathrm{Tx}} + \sum_{\ell=1}^{i} \frac{1}{\lambda_{i,\ell}^2} \right) \right)^{-1}} \\
&= \left( \sum_{\ell=1}^{i} \mu_{\pi(\ell)} \right) \log_2 \left( 1 + \frac{P_{\mathrm{Tx}}}{\sum_{\ell=1}^{i} \frac{1}{\lambda_{i,\ell}^2}} \right).
\end{aligned}
$$

Finally, by using (8) we have

$$\sum_{\ell=1}^{i} \frac{1}{\lambda_{i,\ell}^2} = \sum_{\ell=1}^{i} \left\| \boldsymbol{H}_{\mathrm{co},i}^+ \boldsymbol{e}_\ell \right\|_2^2 = \left\| \boldsymbol{H}_{\mathrm{co},i}^+ \right\|_{\mathrm{F}}^2,$$

which leads to the desired result

$$R_{\mathrm{wsr},i} \geq \left( \sum_{j=1}^{i} \mu_{\pi(j)} \right) \log_2 \left( 1 + \frac{P_{\mathrm{Tx}}}{\left\| \boldsymbol{H}_{\mathrm{co},i}^+ \right\|_{\mathrm{F}}^2} \right).$$

## Appendix B
### Derivation of a Weaker Lower Bound for the Weighted Sum Rate

In this section, we will derive a weaker lower bound for the weighted sum rate $R_{\mathrm{WSR},i}(\pi(1), \dots, \pi(i-1), k, \boldsymbol{b}_1, \dots, \boldsymbol{b}_{i-1}, \hat{\boldsymbol{b}}_i)$ which can be obtained as follows: [see (36) at the bottom of page]. Hence, the receive filters

$$
\begin{aligned}
R_{\mathrm{WSR},i} &\geq \left( \sum_{j=1}^{i-1} \mu_{\pi(j)} + \mu_k \right) \log_2 \left( 1 + \frac{P_{\mathrm{Tx}}}{\mathrm{tr}\left( \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{L}_{i-1}^{-1} \right) + \frac{\hat{\boldsymbol{b}}_i^{\mathrm{H}} \left( \boldsymbol{I}_{r_k} + \boldsymbol{H}_k \boldsymbol{Q}_{i-1} \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{Q}_{i-1}^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \right) \hat{\boldsymbol{b}}_i}{\hat{\boldsymbol{b}}_i^{\mathrm{H}} \boldsymbol{H}_k \left( \boldsymbol{I}_N - \boldsymbol{Q}_{i-1} \boldsymbol{Q}_{i-1}^{\mathrm{H}} \right) \boldsymbol{H}_k^{\mathrm{H}} \hat{\boldsymbol{b}}_i}} \right) \\
&\geq \left( \sum_{j=1}^{i-1} \mu_{\pi(j)} + \mu_k \right) \log_2 \left( 1 + \frac{P_{\mathrm{Tx}}}{\mathrm{tr}\left( \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{L}_{i-1}^{-1} \right) + \frac{1 + \mathrm{tr}\left( \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \right) \max_k \left( \lambda_{\max} \left( \boldsymbol{H}_k, \boldsymbol{H}_k^{\mathrm{H}} \right) \right)}{\hat{\boldsymbol{b}}_i^{\mathrm{H}} \boldsymbol{H}_k \left( \boldsymbol{I}_N - \boldsymbol{Q}_{i-1} \boldsymbol{Q}_{i-1}^{\mathrm{H}} \right) \boldsymbol{H}_k^{\mathrm{H}} \hat{\boldsymbol{b}}_i}} \right). \quad (36)
\end{aligned}
$$

maximizing this lower bound can be determined as described in (24)

$$\tilde{\boldsymbol{b}}_i(k) = \arg\max_{\hat{\boldsymbol{b}}_i} \hat{\boldsymbol{b}}_i^{\mathrm{H}} \boldsymbol{H}_k \left( \boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}} \right) \boldsymbol{H}_k^{\mathrm{H}} \hat{\boldsymbol{b}}_i$$

$$\text{s.t. } \hat{\boldsymbol{b}}_i^{\mathrm{H}} \hat{\boldsymbol{b}}_i = 1.$$

The first inequality is obtained by inserting (18) into (16). The second inequality in (36) is due to the following:

$$\hat{\boldsymbol{b}}_i^{\mathrm{H}} \left( \boldsymbol{I}_{r_k} + \boldsymbol{H}_k \boldsymbol{Q}_{i-1} \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{Q}_{i-1}^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \right) \hat{\boldsymbol{b}}_i$$
$$\leq 1 + \mathrm{tr}\left( \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \right) \max_k \left( \lambda_{\max}\left( \boldsymbol{H}_k, \boldsymbol{H}_k^{\mathrm{H}} \right) \right),$$

which has been derived in [30] and [31]. We will review its derivation here again for the sake of completeness.

$$\hat{\boldsymbol{b}}_i^{\mathrm{H}} \left( \boldsymbol{I}_{r_k} + \boldsymbol{H}_k \boldsymbol{Q}_{i-1} \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{Q}_{i-1}^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \right) \hat{\boldsymbol{b}}_i$$
$$= \hat{\boldsymbol{b}}_i^{\mathrm{H}} \hat{\boldsymbol{b}}_i + \hat{\boldsymbol{b}}_i^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{Q}_{i-1} \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{Q}_{i-1}^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}} \hat{\boldsymbol{b}}_i$$
$$\overset{(a)}{\leq} 1 + \hat{\boldsymbol{b}}_i^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{H}_k^{\mathrm{H}} \hat{\boldsymbol{b}}_i \lambda_{\max}\left( \boldsymbol{Q}_{i-1} \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{Q}_{i-1}^{\mathrm{H}} \right)$$
$$\overset{(b)}{\leq} 1 + \hat{\boldsymbol{b}}_i^{\mathrm{H}} \boldsymbol{H}_k \boldsymbol{H}_k^{\mathrm{H}} \hat{\boldsymbol{b}}_i \mathrm{tr}\left( \boldsymbol{Q}_{i-1} \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{Q}_{i-1}^{\mathrm{H}} \right)$$
$$\overset{(c)}{\leq} 1 + \mathrm{tr}\left( \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \right) \max_k \left( \lambda_{\max}\left( \boldsymbol{H}_k \boldsymbol{H}_k^{\mathrm{H}} \right) \right) \quad (37)$$

which is independent of $\hat{\boldsymbol{b}}_i$ and $k$. (a) stems from the facts that $\hat{\boldsymbol{b}}_i$ is constrained to have norm one throughout the paper and that a quadratic form $\boldsymbol{b}^{\mathrm{H}} \boldsymbol{A} \boldsymbol{b}$ is always smaller or equal to the product of the maximum eigenvalue of $\boldsymbol{A}$ and the norm of $\boldsymbol{b}$, where equality holds if $\boldsymbol{b}$ is a principal eigenvector of $\boldsymbol{A}$. Exploiting the fact that the trace of a positive semi-definite matrix cannot be smaller than the maximum eigenvalue of that matrix leads to (b). Finally (c) is obtained with the same reasoning as (a), the introduction of the max-operator over all users and by using the identities $\mathrm{tr}(\boldsymbol{A}\boldsymbol{B}) = \mathrm{tr}(\boldsymbol{B}\boldsymbol{A})$ and $\boldsymbol{Q}_{i-1}^{\mathrm{H}} \boldsymbol{Q}_{i-1} = \boldsymbol{I}_{i-1}$.

## APPENDIX C
### LOWER AND UPPER BOUNDS FOR A GENERALIZED EIGENVALUE PROBLEM

In this section we will derive a lower and an upper bound for the maximum eigenvalue of a matrix $(\boldsymbol{I}_r + \boldsymbol{D})^{-1}\boldsymbol{C}$, where $\boldsymbol{C} \in \mathbb{C}^{r \times r}$ and $\boldsymbol{D} \in \mathbb{C}^{r \times r}$ are positive semi-definite Hermitian matrices. For (23) we have $\boldsymbol{C} = \boldsymbol{H}_k(\boldsymbol{I}_N - \boldsymbol{Q}_{i-1}\boldsymbol{Q}_{i-1}^{\mathrm{H}})\boldsymbol{H}_k^{\mathrm{H}}$ and $\boldsymbol{D} = \boldsymbol{H}_k \boldsymbol{Q}_{i-1} \boldsymbol{L}_{i-1}^{-1} \boldsymbol{L}_{i-1}^{-\mathrm{H}} \boldsymbol{Q}_{i-1}^{\mathrm{H}} \boldsymbol{H}_k^{\mathrm{H}}$. This maximum eigenvalue can be upper bounded by (e.g. [41], Ch.9, Theorem H.1.a)

$$\lambda_{\max}\left( (\boldsymbol{I}_r + \boldsymbol{D})^{-1}\boldsymbol{C} \right) \leq \lambda_{\max}\left( (\boldsymbol{I}_r + \boldsymbol{D})^{-1} \right) \lambda_{\max}(\boldsymbol{C}).$$

As the minimum eigenvalue of the matrix $\boldsymbol{I}_r + \boldsymbol{D}$ is greater than or equal to one, when $\boldsymbol{D}$ is positive semi-definite, the maximum eigenvalue of its inverse is equal to or smaller than one. Furthermore the maximum eigenvalue of a positive semi-definite matrix $\boldsymbol{C}$ is smaller than or equal to the trace of that matrix which leads to the upper bound

$$\lambda_{\max}\left( (\boldsymbol{I}_r + \boldsymbol{D})^{-1}\boldsymbol{C} \right) \leq \mathrm{tr}(\boldsymbol{C}).$$

For the derivation of the lower bound we use the lower bound for the maximum eigenvalue from [37], Ch. 2.3

$$\frac{\mathrm{tr}\left( (\boldsymbol{I}_r + \boldsymbol{D})^{-1}\boldsymbol{C} \right)}{r} \leq \lambda_{\max}\left( (\boldsymbol{I}_r + \boldsymbol{D})^{-1}\boldsymbol{C} \right),$$

where equality holds, if all eigenvalues are identical. Denoting the eigenvalue decompositions of the matrices $\boldsymbol{C}$ and $(\boldsymbol{I}_r + \boldsymbol{D})^{-1}$ as

$$\boldsymbol{C} = \sum_{i=1}^{\mathrm{rank}(\boldsymbol{C})} \lambda_i \boldsymbol{u}_i \boldsymbol{u}_i^{\mathrm{H}}, \quad (\boldsymbol{I}_r + \boldsymbol{D})^{-1} = \sum_{j=1}^{r} \frac{1}{1 + \rho_j} \boldsymbol{v}_j \boldsymbol{v}_j^{\mathrm{H}},$$

where the $\lambda_j$'s and the $\rho_j$'s correspond to the eigenvalues of the matrix $\boldsymbol{C}$ and the matrix $\boldsymbol{D}$, respectively, we obtain

$$\mathrm{tr}\left( (\boldsymbol{I}_r + \boldsymbol{D})^{-1}\boldsymbol{C} \right) = \sum_{i=1}^{\mathrm{rank}(\boldsymbol{C})} \sum_{j=1}^{r} \frac{\lambda_i}{1 + \rho_j} \mathrm{tr}\left( \boldsymbol{u}_i \boldsymbol{u}_i^{\mathrm{H}} \boldsymbol{v}_j \boldsymbol{v}_j^{\mathrm{H}} \right)$$
$$= \sum_{i=1}^{\mathrm{rank}(\boldsymbol{C})} \sum_{j=1}^{r} \frac{\lambda_i}{1 + \rho_j} \left| \boldsymbol{v}_j^{\mathrm{H}} \boldsymbol{u}_i \right|^2 \quad (38)$$

As $\boldsymbol{D}$ is positive semi-definite, all its eigenvalues $\rho_j$ can be upper bounded by its trace, i.e., $\rho_j \leq \mathrm{tr}(\boldsymbol{D})$. Hence, we can lower bound the expression (38) as

$$\sum_{i=1}^{\mathrm{rank}(\boldsymbol{C})} \sum_{j=1}^{r} \frac{\lambda_i}{1 + \rho_j} \left| \boldsymbol{v}_j^{\mathrm{H}} \boldsymbol{u}_i \right|^2 \geq \sum_{i=1}^{\mathrm{rank}(\boldsymbol{C})} \sum_{j=1}^{r} \frac{\lambda_i}{1 + \mathrm{tr}(\boldsymbol{D})} \left| \boldsymbol{v}_j^{\mathrm{H}} \boldsymbol{u}_i \right|^2$$
$$= \sum_{i=1}^{\mathrm{rank}(\boldsymbol{C})} \frac{\lambda_i}{1 + \mathrm{tr}(\boldsymbol{D})} \sum_{j=1}^{r} \left| \boldsymbol{v}_j^{\mathrm{H}} \boldsymbol{u}_i \right|^2.$$

As the vectors $\boldsymbol{v}_j$ form an orthonormal basis, we obtain

$$\sum_{i=1}^{\mathrm{rank}(\boldsymbol{C})} \frac{\lambda_i}{1 + \mathrm{tr}(\boldsymbol{D})} \sum_{j=1}^{r} \left| \boldsymbol{v}_j^{\mathrm{H}} \boldsymbol{u}_i \right|^2 = \sum_{i=1}^{\mathrm{rank}(\boldsymbol{C})} \frac{\lambda_i}{1 + \mathrm{tr}(\boldsymbol{D})}$$
$$= \frac{\mathrm{tr}(\boldsymbol{C})}{1 + \mathrm{tr}(\boldsymbol{D})}, \quad (39)$$

which is the desired result.

### REFERENCES

[1] 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, Sep. 2006.

[2] H. Viswanathan, S. Venkatesan, and H. Huang, "Downlink capacity evaluation of cellular networks with known-interference cancellation," *IEEE J. Sel. Areas Commun.*, vol. 21, no. 6, pp. 802–811, Jun. 2003.

[3] R. Hunger, D. Schmidt, M. Joham, and W. Utschick, "A general covariance-based optimization framework using orthogonal projections," in *Proc. IEEE Conf. Signal Process. Adv. Wireless Commun. (SPAWC)*, Jul. 2008, pp. 76–80.

[4] R. Böhnke and K.-D. Kammeyer, "Weighted sum rate maximization for MIMO-OFDM systems with linear and dirty paper precoding," in *Proc. 7th Int. ITG Conf. Source and Channel Coding (SCC 08)*, Ulm, Germany, Jan. 2008, 5 pages.

[5] M. Kobayashi and G. Caire, "An iterative water-filling algorithm for maximum weighted sum-rate of Gaussian MIMO-BC," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 8, pp. 1640–1646, Aug. 2006.

[6] J. Lee and N. Jindal, "Symmetric capacity of MIMO downlink channels," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jul. 2006, pp. 1031–1035.

[7] G. Wunder and T. Michel, "Minimum rates scheduling for MIMO-OFDM broadcast channels," in *Proc. 9th IEEE Int. Symp. Spread Spectrum Techn. Appl. (ISSSTA 2006)*, Aug. 2006, pp. 510–514.

[8] C. Fung, W. Yu, and T. Lim, "Multiantenna downlink precoding with individual rate constraints: Power minimization and user ordering," in *Proc. Int. Conf. Commun. Syst.*, Sep. 2004, pp. 45–49.

[9] M. Costa, "Writing on dirty paper," *IEEE Trans. Inf. Theory*, vol. 29, no. 3, pp. 439–441, May 1983.

[10] D. Schmidt, M. Joham, and W. Utschick, "Minimum mean square error vector precoding," *Eur. Trans. Telecommun.*, vol. 19, pp. 219–231, 2007.

[11] S. ten Brink and U. Erez, "A close-to-capacity dirty paper coding scheme," in *Proc. Int. Symp. Inf. Theory (ISIT)*, Jul. 2004, p. 533.

[12] M. Tomlinson, "New automatic equalizer employing modulo arithmetic," *Electron. Lett.*, vol. 7, pp. 138–139, March 1971.

[13] H. Harashima and H. Miyakawa, "Matched-transmission technique for channels with intersymbol interference," *IEEE Trans. Commun.*, vol. 20, no. 4, pp. 774–780, Aug. 1972.

[14] W. Yu, D. P. Varodayan, and J. M. Cioffi, "Trellis and convolutional precoding for transmitter-based interference presubtraction," *IEEE Trans. Commun.*, vol. 53, no. 7, pp. 1120–1230, Jul. 2005.

[15] R. Hunger, D. A. Schmidt, and M. Joham, "A combinatorial approach to maximizing the sum rate in the MIMO BC with linear filtering," in *Proc. 42nd Asilomar Conf. Signals, Syst., and Comput. 2008*, Oct. 2008, p. 5 pages.

[16] S. Shi, M. Schubert, and H. Boche, "Rate optimization for multiuser MIMO systems with linear precoding," *IEEE Trans. Signal Process.*, vol. 56, no. 8, pp. 4020–4030, Aug. 2008.

[17] E. Jorswieck and E. Larsson, "Linear precoding in multiple antenna broadcast channels: Efficient computation of the achievable rate region," in *Proc. IEEE/ITG Workshop on Smart Antennas (WSA)*, Feb. 2008, pp. 21–28.

[18] S. Shi, M. Schubert, and H. Boche, "Downlink MMSE transceiver optimization for multiuser MIMO systems: Duality and sum-MSE minimization," *IEEE Trans. Signal Process.*, vol. 55, no. 11, pp. 5436–5446, 2007.

[19] R. Hunger and M. Joham, "A general rate duality of the MIMO multiple access channel and the MIMO broadcast channel," in *Proc. IEEE Global Telecommun. Conf. (GLOBECOM)*, Dec. 2008, pp. 1–5.

[20] M. Codreanu, A. Tölli, M. Juntti, and M. Latva-aho, "Joint design of Tx-Rx beamformers in MIMO downlink channel," *IEEE Trans. Signal Process.*, vol. 55, no. 9, pp. 4639–4655, Sep. 2007.

[21] Q. Spencer, A. Swindlehurst, and M. Haardt, "Zero-forcing methods for downlink spatial multiplexing in multiuser MIMO channels," *IEEE Trans. Signal Process.*, vol. 52, no. 2, pp. 461–471, Feb. 2004.

[22] G. Caire and S. Shamai, "On the achievable throughput of multiantenna Gaussian broadcast channel," *IEEE Trans. Inf. Theory*, vol. 49, no. 7, pp. 1691–1706, Jul. 2003.

[23] M. Fuchs, G. D. Galdo, and M. Haardt, "Low-complexity space-time-frequency scheduling for MIMO systems with SDMA," *IEEE Trans. Veh. Technol.*, vol. 56, no. 5, pp. 2775–2784, Sep. 2007.

[24] Z. Shen, R. Chen, J. Andrews, R. Heath, and B. Evans, "Low complexity user selection algorithms for multiuser MIMO systems with block diagonalization," *IEEE Trans. Signal Process.*, vol. 54, no. 9, pp. 3658–3663, Sep. 2006.

[25] R. Chen, Z. Shen, J. Andrews, and R. Heath, "Multimode transmission for multiuser MIMO systems with block diagonalization," *IEEE Trans. Signal Process.*, vol. 56, no. 7, pp. 3294–3302, Jul. 2008.

[26] G. Dimić and N. Sidoropoulos, "On downlink beamforming with greedy user selection," *IEEE Trans. Signal Process.*, vol. 53, no. 10, pp. 3857–3868, Oct. 2005.

[27] J. Wang, D. Love, and M. Zoltowski, "User selection with zero-forcing beamforming achieves the asymptotically optimal sum rate," *IEEE Trans. Signal Process.*, vol. 56, no. 8, pp. 3713–3726, Aug. 2008.

[28] T. Yoo and A. Goldsmith, "On the optimality of multiantenna broadcast scheduling using zero-forcing beamforming," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 3, pp. 528–541, Mar. 2006.

[29] F. Boccardi and H. Huang, "A near-optimum technique using linear precoding for the MIMO broadcast channel," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Apr. 2007, vol. 3, pp. III-17–III-20.

[30] C. Guthy, W. Utschick, G. Dietl, and P. Tejera, "Efficient linear successive allocation for the MIMO broadcast channel," in *Proc. 42nd Asilomar Conf. Signals, Syst., Comput.*, Oct. 2008.

[31] C. Guthy, W. Utschick, and G. Dietl, "Low complexity linear zero-forcing for the MIMO broadcast channel," *IEEE J. Sel. Topics in Signal Process., Special Issue on Managing Complexity in Multiuser MIMO Syst.*, Dec. 2009.

[32] Y. Hara, L. Brunel, and K. Oshima, "Spatial scheduling with interference cancellation in multiuser MIMO systems," *IEEE Trans. Veh. Technol.*, vol. 57, no. 2, pp. 893–905, Mar. 2008.

[33] J. Wang, D. Love, and M. Zoltowski, "User selection for the MIMO broadcast channel with a fairness constraint," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2007, pp. III-9–III-12.

[34] A. Wiesel, Y. Eldar, and S. Shamai, "Zero-forcing precoding and generalized inverses," *IEEE Trans. Inf. Theory*, vol. 56, no. 9, pp. 4409–4418, Sep. 2008.

[35] M. Bazaraa, H. Sherali, and C. Shetty, *Nonlinear Programming-Theory and Applications*, 3rd ed. New York: Wiley Interscience, 2006.

[36] C. Swannack, E. Uysal-Biyikoglu, and G. Wornell, "Low complexity multiuser scheduling for maximizing throughput in the MIMO broadcast channel," in *Proc. 42nd Ann. Allerton Conf. Commun., Control, Comput.*, 2004, pp. 440–449.

[37] G. H. Golub and C. F. van Loan, *Matrix Computations*. Baltimore, MD: The John Hopkins Univ. Press, 1989.

[38] C. Guthy, W. Utschick, J. A. Nossek, G. Dietl, and G. Bauch, "Rate-invariant user preselection for complexity reduction in multiuser MIMO systems," in *Proc. IEEE Veh. Technol. Conf. (VTC)*, Sep. 2008.

[39] P. Tejera, W. Utschick, G. Bauch, and J. Nossek, "Efficient implementation of successive encoding schemes for the MIMO OFDM broadcast channel," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Istanbul, Jun. 2006.

[40] C.-B. Chae, D. Mazzarese, T. Inoue, and R. Heath, "Coordinated beamforming for the multiuser MIMO broadcast channel with limited feedforward," *IEEE Trans. Signal Process.*, vol. 56, no. 12, pp. 6044–6056, Dec. 2008.

[41] A. W. Marshall and I. Olkin, *Inequalities: Theory of Majorization and Its Applications*, ser. Math. Sci. Eng., R. Bellman, Ed. New York: Academic, 1979, vol. 143.

**Christian Guthy** (S'10) was born in Munich, Germany, in 1979. He received the B.Sc. and Dipl.-Ing. degrees (the latter with *summa cum laude*) in electrical engineering from the Technische Universität München (TUM), Munich, in 2004 and 2005, respectively.

Since 2005, he has been with the Associate Institute for Signal Processing, TUM, where he is currently working toward the Ph.D degree. His main research interests include design and analysis of low complexity signal processing algorithms for next generation wireless communication systems with focus on multiantenna and multicarrier systems.

**Wolfgang Utschick** (SM'06) was born on May 6, 1964. He completed several industrial education programs before he received the diploma and doctoral degrees, both with honors, in electrical engineering from the Technische Universität München, Germany (TUM), in 1993 and 1998, respectively. In this period he held a scholarship of the Bavarian Ministry of Education for exceptional students.

From 1998 to 2002, he codirected the Signal Processing Group of the Institute of Circuit Theory and Signal Processing, TUM. Since 2000, he has been consulting in 3 GPP standardization in the field of multielement antenna systems. In 2002, he was appointed Professor at the TUM where he is Head of the Fachgebiet Methoden der Signalverarbeitung. He teaches courses on signal processing, stochastic processes, and optimization theory in the field of digital communications.

Dr. Utschick was awarded in 2006 for his excellent teaching records at TUM, and in 2007, he received the ITG Award of the German Society for Information Technology (ITG). He is a senior member of the German VDE/ITG where he has been appointed in the Expert Committee for Information and System Theory. He is currently also serving as an Associate Editor for the IEEE TRANSACTIONS ON SIGNAL PROCESSING.

**Raphael Hunger** (S'06) was born in Illertissen, Germany, 1979. He studied electrical engineering at the Technische Universität München (TUM), Germany, from 1999 until 2005 and, in 2004, at the RWTH Aachen, Germany. He received the Dipl.-Ing. and Master of Science degrees in electrical engineering from the TUM in 2004 and 2005, respectively.

Since 2005, he has been working toward the doctorate degree at the Associate Institute for Signal Processing, TUM. His research interests focus on the joint optimization of transmitters and receivers in multiuser MIMO communications.

**Michael Joham** (S'99–M'05) was born in Kufstein, Austria, 1974. He received the Dipl.-Ing. and Dr.-Ing. degrees (both *summa cum laude*) in electrical engineering from the Technische Universität München (TUM), Germany, in 1999 and 2004, respectively.

He was with the Institute for Circuit Theory and Signal Processing, TUM, from 1999 to 2004. Since 2004, he has been with the Associate Institute for Signal Processing, TUM, where he is currently a Senior Researcher. During summers 1998 and 2000, he visited Purdue University, West Lafayette, IN. In spring 2008, he was a Guest Lecturer with the University of the German Federal Armed Forces, Munich, and a guest professor with the University of A Coruna, Spain. In winter 2009, he was a Guest Lecturer with the University of Hanover, Germany. His current research interests are precoding in mobile and satellite communications, limited rate feedback, MIMO communications, and robust signal processing.

Dr. Joham received the VDE Preis for his diploma thesis in 1999 and the Texas-Instruments-Preis for his dissertation in 2004. In 2007, he was a corecipient of the Best Paper Award at the International ITG/IEEE Workshop on Smart Antennas in Vienna.