1

# Traffic Aware Resource Allocation Schemes for Multi-Cell MIMO-OFDM Systems

Ganesh Venkatraman Student Member, IEEE, Antti Tölli Senior Member, IEEE, Markku Juntti Senior Member, IEEE, and Le-Nam Tran Member, IEEE

Abstract—We consider a downlink multi-cell multiple-input multiple-output (MIMO) interference broadcast channel (IBC) using orthogonal frequency division multiplexing (OFDM) with multiple users contending for space-frequency resources in a given scheduling instant. The problem is to design precoders efficiently to minimize the number of backlogged packets queuing in the coordinating base stations (BSs). Conventionally, the queue weighted sum rate maximization (Q-WSRM) formulation with the number of backlogged packets as the corresponding weights is used to design the precoders. In contrast, we propose joint spacefrequency resource allocation (JSFRA) formulation, in which the precoders are designed jointly across the space-frequency resources for all users by minimizing the total number of backlogged packets in each transmission instant, thereby performing user scheduling implicitly. Since the problem is nonconvex, we use the combination of successive convex approximation (SCA) and alternating optimization (AO) to handle nonconvex constraints in the JSFRA formulation. In the first method, we approximate the signal-to-interference-plus-noise ratio (SINR) by convex relaxations, while in the second approach, the equivalence between the SINR and the mean squared error (MSE) is exploited. We then discuss the distributed approaches for the centralized algorithms using primal decomposition and alternating directions method of multipliers (ADMM). Finally, we propose a practical iterative precoder design by solving the Karush-Kuhn-Tucker expressions for the MSE reformulation that requires minimal information exchange for each update. Numerical results are used to compare the proposed algorithms to the existing solutions.

*Index Terms*—Convex approximations, MIMO-IBC, MIMO-OFDM, precoder design, SCA, WSRM.

#### I. INTRODUCTION

In a network with multiple base stations (BSs) serving multiple users, the main driving factor for the transmission is the packets waiting at each BS corresponding to the different users present in the network. We consider the problem of transmit precoder design over the space-frequency resources provided by the multiple-input multiple-output (MIMO) orthogonal frequency division multiplexing (OFDM) framework in the downlink interference broadcast channel (IBC) to minimize the number of queued packets at each BS. Since the available resources are shared by multiple users associated with different BSs, it can be viewed as a resource allocation problem.

This work has been supported by the Finnish Funding Agency for Innovation (Tekes), Nokia Networks, Xilinx, Elektrobit and the Academy of Finland. Part of this work is presented in ICASSP 2014.

- G. Venkatraman, A. Tölli and M. Juntti are with Centre for Wireless Communications (CWC), Department of Communications Engineering (DCE), University of Oulu, Oulu, FI-90014, (e-mail: {ganesh.venkatraman, antti.tolli, markku.juntti}@ee.oulu.fi)
- L.-N. Tran was with Centre for Wireless Communications (CWC), Oulu, FI-90014. He is now with the Department of Electronic Engineering, Maynooth University, Maynooth, Co. Kildare, Ireland, (e-mail:ltran@eeng.nuim.ie)

In general, the resource allocation problems such as admission control ones can be formulated by assigning a binary variable for each user to indicate the presence or the absence in a particular resource [1]. Alternatively, linear transmit precoders, which are complex vectors, can be implicitly modeled as decision variables, thereby avoiding the use of binary decision variables. After the decision stage, the non-zero precoders are used to determine the transmission rates of users on a space-frequency resource. A zero transmit precoder indicates the absence of the user on a given resource. In this way, the soft decisions are used in the optimization problem and the hard decisions are made after the algorithm convergences.

The queue minimizing precoder designs are closely related to the weighted sum rate maximization (WSRM) problem with additional rate constraints determined by the number of backlogged packets for each user in the system. The topics on MIMO IBC precoder design have been studied extensively with different performance criteria in the literature. Due to the nonconvex nature of the MIMO IBC precoder design problems, the successive convex approximation (SCA) approach has become a powerful tool to deal with these problems. For example, in [2], the nonconvex part of the objective has been linearized around an operating point in order to solve the WSRM problem in an iterative manner. Similar approach of solving the WSRM problem by using arithmetic-geometric inequality has been proposed in [3].

The relation between the achievable sum rate and the mean squared error (MSE) of the received symbol by using fixed minimum mean squared error (MMSE) receivers can be used to solve the WSRM problem [4]. In [5], [6], the WSRM problem is reformulated via MSE, casting the problem as a convex one for fixed linearization coefficients. In this way, the original problem is expressed in terms of the MSE weight, precoders, and decoders. Then the problem is solved using an alternating optimization method, i.e., finding a subset of variables while the other variables are fixed. The MSE reformulation for the WSRM problem has also been studied in [7] by using SCA to solve the problem in an iterative manner. Moreover, distributed precoder design with quality of service (QoS) requirements as additional rate constraints are studied for the MSE reformulated WSRM problem in [8], [9].

The problem of precoder design for the MIMO IBC system can be solved either by using a centralized controller or by using decentralized algorithms, where each BS handles the corresponding subproblem independently with the limited information exchange with other BSs via back-haul. The distributed approaches are usually based on the primal, or

dual decompositions or the alternating directions method of multipliers (ADMM), which has been discussed in [10], [11]. In the primal decomposition, the so-called coupling interference variables are fixed for the subproblem at each BS to find the optimal precoders. The fixed interference is then updated by using the subgradient method as discussed in [12]. The dual and the ADMM approaches control the distributed subproblems by fixing the 'interference price' for each BS as detailed in [13].

By adjusting the weights in the WSRM objective, we can find an arbitrary rate-tuple in the rate region that maximizes suitable objective measures. For example, if the weight of each user is set to be inversely proportional to its average data rate, the corresponding problem guarantees fairness on the average among the users. To reduce the number of backlogged packets, we can assign weights based on the current queue size of the users. Specifically, the queue states can be incorporated in the WSRM objective  $\sum_k w_k R_k$  by replacing the weight  $w_k$  with the corresponding queue state  $Q_k$  or its function, which is the outcome of minimizing the Lyapunov drift between the current and the future queue states [14], [15], where  $R_k$  denotes the achievable data rate of user k. In the backpressure algorithm, the differential queues between the source and the receiver nodes are used to scale the transmission rate [16].

Earlier studies on the queue minimization problem are summarized in the survey papers [17], [18]. In particular, the problem of power allocation to minimize the number of backlogged packets is considered in [19] using geometric programming. Since the problem addressed in [19] assumed single antenna transmitters and receivers, the queue minimizing problem reduces to the optimal power allocation problem. In the context of wireless networks, the *backpressure algorithm* mentioned above is extended in [20] by formulating the corresponding user queues as the weights in the WSRM problem. Recently, the precoder design for the video transmission over MIMO system is considered in [21]. In this design, the MU-MIMO precoders are designed by the MSE reformulation as in [5] with the higher layer performance objective such as playback interruptions and buffer overflow probabilities.

Main Contributions: In this paper, we design precoders jointly over space-frequency resources to reduce the number of backlogged packets waiting at each BSs. The proposed formulation also limits the allocations beyond the number of backlogged packets without explicit rate constraints. Initially, we propose a centralized joint space-frequency resource allocation (JSFRA) formulation, which is solved by two iterative algorithms based on the combination of SCA and alternating optimization (AO) due to the nonconvex nature of the problem. The proposed algorithms solve a sequence of convex problems obtained by fixing a subset of optimization variables or by approximating the nonconvex constraints by the convex ones. The first approach is performed by directly relaxing the signalto-interference-plus-noise ratio (SINR) expression, while in the second method, the equivalence between the MSE and the SINR is exploited. We then discuss the distributed implementation of the JSFRA methods using primal decomposition and the ADMM. Finally, we also propose a more practical iterative precoder design by directly solving the Karush-Kuhn-Tucker (KKT) system of equations for the MSE reformulation that is numerically shown to require minimal information exchange for each update. Note that the joint space-frequency channel matrix can be formed by stacking the channel of each subchannel in a block-diagonal form for all users.

The rest of the paper is organized as follows. In Section II, we introduce the system model and the problem formulation for the queue minimizing precoder design. The existing and the proposed centralized precoder designs are presented in Section III. The distributed solutions are provided in Section IV followed by the simulation results in Section V. Conclusions are drawn in Section VI.

### II. SYSTEM MODEL AND PROBLEM FORMULATION

#### A. System Model

We consider a downlink MIMO IBC scenario in an OFDM framework with N sub-channels and  $N_B$  BSs each equipped with  $N_T$  transmit antennas, serving in total K users each with  $N_R$  receive antennas. The set of users associated with BS b is denoted by  $\mathcal{U}_b$  and the set  $\mathcal{U}$  represents all users in the system, i.e.,  $\mathcal{U} = \bigcup_{b \in \mathcal{B}} \mathcal{U}_b$ , where  $\mathcal{B}$  is the set of indices of all coordinating BSs. Data for user k is transmitted from only one BS which is denoted by  $b_k \in \mathcal{B}$ . Let  $\mathcal{N} = \{1, 2, \ldots, N\}$  be the set of all sub-channel indices available in the system.

We adopt linear transmit beamforming technique at BSs. Specifically, the data symbols  $d_{l,k,n}$  for user k on the  $l^{\rm th}$  spatial stream over sub-channel n is multiplied with beamformer  $\mathbf{m}_{l,k,n} \in \mathbb{C}^{N_T \times 1}$  before being transmitted. In order to detect multiple spatial streams at the user terminal, receive beamforming vector  $\mathbf{w}_{l,k,n}$  is employed for each user. Consequently, the received data estimate corresponding to the  $l^{\rm th}$  spatial stream over sub-channel n at user k is given by

$$\hat{d}_{l,k,n} = \mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b_k,k,n} \, \mathbf{m}_{l,k,n} d_{l,k,n} + \mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{n}_{k,n} + \mathbf{w}_{l,k,n}^{\mathrm{H}} \sum_{i \in \mathcal{U} \setminus \{k\}} \mathbf{H}_{b_i,k,n} \sum_{j=1}^{L} \mathbf{m}_{j,i,n} d_{j,i,n} \quad (1)$$

where  $\mathbf{H}_{b,k,n} \in \mathbb{C}^{N_R \times N_T}$  is the channel between BS b and user k on sub-channel n, and  $\mathbf{n}_{k,n} \sim \mathcal{CN}(0,N_0)$  is the additive noise vector for user k on the  $n^{\mathrm{th}}$  sub-channel and  $l^{\mathrm{th}}$  spatial stream. In (1),  $L = \mathrm{rank}(\mathbf{H}_{b,k,n}) = \min(N_T, N_R)$  is the maximum number of spatial streams. Assuming independent detection of data streams, we can write the signal-to-interference-plus-noise ratio (SINR) as

$$\gamma_{l,k,n} = \frac{\left| \mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b_k,k,n} \mathbf{m}_{l,k,n} \right|^2}{\widetilde{N}_0 + \sum_{(j,i) \neq (l,k)} |\mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b_i,k,n} \mathbf{m}_{j,i,n}|^2}$$
(2)

where  $\widetilde{N}_0 = N_0 \operatorname{tr}(\mathbf{w}_{l,k,n} \mathbf{w}_{l,k,n}^{\mathrm{H}})$  denotes the equivalent noise variance. To reduce the overhead involved in feeding back the user channels, we consider a time division duplexing (TDD) system in which BSs can estimate the downlink channels from the uplink pilots using channel reciprocity.

 $<sup>^1</sup>$ It can be easily extended for user specific streams  $L_k$  instead of using common L streams for all users. L streams are initialized but after solving the problem, only  $L_{k,n} \leq L$  non-zero data streams are transmitted.

Let  $Q_k$  be the number of backlogged packets destined for user k at a given scheduling instant. The queue dynamics of user k are modeled using the Poisson arrival process with the average number of packet arrivals of  $A_k = \mathbf{E}_i\{\lambda_k\}$  packets/bits, where  $\lambda_k(i) \sim \operatorname{Pois}(A_k)$  represents the instantaneous number of packets arriving for user k at the  $i^{\text{th}}$  time instant. The total number of queued packets at the  $(i+1)^{\text{th}}$  instant for user k, denoted as  $Q_k(i+1)$ , is given by

$$Q_k(i+1) = [Q_k(i) - t_k(i)]^+ + \lambda_k(i)$$
 (3)

where  $[x]^+ \equiv \max\{x,0\}$  and  $t_k$  denotes the number of transmitted packets or bits for user k. At the  $i^{\text{th}}$  instant, the transmission rate of user k is given by

$$t_k(i) = \sum_{n=1}^{N} \sum_{l=1}^{L} t_{l,k,n}(i)$$
 (4)

where  $t_{l,k,n}$  denotes the number of transmitted packets or bits over the  $l^{\rm th}$  spatial stream on the  $n^{\rm th}$  sub-channel. The maximum rate achieved over the space-frequency resource (l,n) is given by  $t_{l,k,n} \leq \log_2(1+\gamma_{l,k,n})$  for the SINR  $\gamma_{l,k,n}$ . The units of  $t_k$  and  $Q_k$  are in bits defined per channel use.

#### B. Problem Formulation

To minimize the total number of backlogged packets, we consider minimizing the weighted  $\ell_q$ -norm of the queue deviation objective as

$$v_k = Q_k - t_k = Q_k - \sum_{n=1}^{N} \sum_{l=1}^{L} \log_2(1 + \gamma_{l,k,n})$$
 (5)

where  $\gamma_{l,k,n}$  is given by (2) and the optimization variables are the transmit precoders  $\mathbf{m}_{l,k,n}$  and the receivers  $\mathbf{w}_{l,k,n}$ .

Explicitly, the objective of the problem considered is given as  $\sum_{k\in\mathcal{U}} a_k |v_k|^q$ . Thus the formulation becomes

$$\underset{\mathbf{m}_{l,k,n},\mathbf{w}_{l,k,n}}{\text{minimize}} \quad \|\tilde{\mathbf{v}}\|_{q} \tag{6a}$$

subject to 
$$\sum_{n=1}^{N} \sum_{k \in \mathcal{U}_b} \sum_{l=1}^{L} \operatorname{tr}\left(\mathbf{m}_{l,k,n} \mathbf{m}_{l,k,n}^{\mathrm{H}}\right) \leq P_{\max} \forall b \ \ (6b)$$

where  $\tilde{v}_k \triangleq a_k^{1/q} v_k$  is the element of vector  $\tilde{\mathbf{v}}$ , and  $a_k$  is the weighting factor, which is used to alter the user priority based on the QoS constraints such as packet delay requirements and packet waiting time, since they are proportional to the corresponding number of backlogged packets. The BS specific power constraint for all sub-channels is considered in (6b).

For practical reasons, we impose a constraint on the maximum number of transmitted bits for the user k, since it is limited by the total number of backlogged packets available at the transmitter. As a result, the number of backlogged packets  $v_k$  for user k remaining in the system is given by

$$v_k = Q_k - \sum_{n=1}^{N} \sum_{l=1}^{L} \log_2(1 + \gamma_{l,k,n}) \ge 0.$$
 (7)

The above positivity constraint need to be satisfied by  $v_k$  to avoid the excessive allocation of the resources.

Before proceeding further, we show that the constraint in (7) is handled implicitly by the definition of norm  $\ell_q$  in the objective of (6). Suppose that  $t_k > Q_k$  for certain k at the optimum, i.e.,  $-v_k = t_k - Q_k > 0$ . Then there exists  $\delta_k > 0$  such that  $-v_k' = t_k' - Q_k < -v_k$  where  $t_k' = t_k - \delta_k$ . Since  $\|\tilde{\mathbf{v}}\|_q = \||\tilde{\mathbf{v}}\|_q = \||-\tilde{\mathbf{v}}\|_q$ , this means that the newly created vector  $\mathbf{t}'$  achieves a smaller objective which contradicts with the fact that the optimal solution has been obtained. The choice of the norm  $\ell_q$  used in the objective function [17], [19] alters the priorities for the queue deviation function as follows.

- $\ell_1$  results in greedy allocation *i.e.*, emptying the queue of users with good channel conditions before considering the users with worse channel conditions. As a special case, it is easy to see that (6) reduces to the WSRM problem when the queue size is large enough for all users.
- l<sub>2</sub> prioritizes users with a higher number of queued packets before considering the users with a smaller number of backlogged packets. For example, it could be more ideal for the delay limited scenario when the packet arrival rates of the users are similar, since the number of backlogged packets is proportional to the delay in the transmission following the Little's law [15].
- $\ell_{\infty}$  minimizes the maximum number of queued packets among users with the current transmission, thereby providing queue fairness by allocating the resources proportional to the number of backlogged packets.

#### III. PROPOSED QUEUE MINIMIZING PRECODER DESIGNS

In general, the precoder design for the MIMO OFDM problem is difficult due to its nonconvex nature. In addition, the objective of minimizing the number of the queued packets over space-frequency dimensions adds further complexity. Since the scheduling of users in each sub-channel attained by allocating zero transmit power over certain sub-channels, our solutions perform joint precoder design and user scheduling. Before discussing the proposed solutions, we consider the existing algorithm to minimize the number of backlogged packets with additional constraints required by the problem.

## A. Queue Weighted Sum Rate Maximization (Q-WSRM) Formulation

The queue minimizing algorithms have been extensively discussed in the networking literature to provide congestion-free routing between any two nodes in the network. One such algorithm is the *backpressure algorithm* [14]–[16]. It determines an optimal control policy in the form of rate or resource allocation for the nodes in the network by considering the differential backlogged packets between the source and the destination nodes. Even though the algorithm is primarily designed for the wired infrastructure, it can be extended to the wireless networks by designing the user rate variable  $t_k$  in accordance to the wireless network.

The queue weighted sum rate maximization (Q-WSRM) formulation extends the *backpressure algorithm* to the downlink MIMO-OFDM framework, in which the multiple BSs act as

<sup>&</sup>lt;sup>2</sup>The unit can either be packets or bits as long as the arrival and the transmission units are similar.

<sup>&</sup>lt;sup>3</sup>The upper bound can be achieved by using Gaussian signaling.

the source nodes and the user terminals as the receiver nodes. The control policy in the form of transmit precoders aims at minimizing the number of queued packets waiting in the BSs. In order to find the optimal strategy, we resort to the Lyapunov theory, which is predominantly used in the control theory to achieve system stability. Since at each time slot, the system is described by the channel conditions and the number of backlogged packets of each user, the Lyapunov function is used to provide a scalar measure, which grows large when the system moves toward the undesirable state. By following [15], the scalar measure for the queue stability is given by

$$L\left[\mathbf{Q}(i)\right] = \frac{1}{2} \sum_{k \in \mathcal{U}} Q_k^2(i) \tag{8}$$

where  $\mathbf{Q}(i) = [Q_1(i), Q_2(i), \dots, Q_K(i)]^T$  and  $\frac{1}{2}$  is used for the convenience. It provides a scalar measure of congestion present in the system [15, Ch. 3].

To minimize the total number of backlogged packets for time instant i, the optimal transmission rate of all users is obtained by minimizing the Lyapunov drift expressed as

$$L[\mathbf{Q}(i+1)] - L[\mathbf{Q}(i)] = \frac{1}{2} \left[ \sum_{k \in \mathcal{U}} \left( [Q_k(i) - t_k(i)]^+ + \lambda_k(i) \right)^2 - Q_k^2(i) \right].$$
(9)

In order to eliminate the nonlinear operator  $[x]^+$ , we bound the expression in (9) as

$$\leq \sum_{k \in \mathcal{U}} \frac{\lambda_k^2(i) + t_k^2(i)}{2} + \sum_{k \in \mathcal{U}} Q_k(i) \left\{ \lambda_k(i) - t_k(i) \right\} \tag{10}$$

by using the following inequality

$$[\max(Q-t,0)+\lambda)]^2 \le Q^2+t^2+\lambda^2+2Q(\lambda-t).$$
 (11)

The total number of backlogged packets at any given instant i is reduced by minimizing the conditional expectation of the Lyapunov drift expression (10) given the current number of queued packets Q(i) waiting in the system. The expectation is taken over all possible arrival and transmission rates of the users to obtain the optimal rate allocation strategy.

Now, the conditional Lyapunov drift, denoted by  $\Delta(Q(i))$ , is given by the infimum over the transmission rate as

$$\inf_{\mathbf{t}} \quad \mathbb{E}_{\boldsymbol{\lambda}, \mathbf{t}} \left\{ \operatorname{L} \left[ \mathbf{Q}(i+1) \right] - \operatorname{L} \left[ \mathbf{Q}(i) \right] | \mathbf{Q}(i) \right\} \tag{12a}$$

$$\leq \underbrace{\mathbb{E}_{\boldsymbol{\lambda}, \mathbf{t}} \left\{ \sum_{k \in \mathcal{U}} \frac{\lambda_k^2(i) + t_k^2(i)}{2} | \mathbf{Q}(i) \right\}}_{\leq B} + \sum_{k \in \mathcal{U}} Q_k(i) A_k(i)$$

$$- \mathbb{E}_{\boldsymbol{\lambda}, \mathbf{t}} \left\{ \sum_{k \in \mathcal{U}} Q_k(i) t_k(i) | \mathbf{Q}(i) \right\}, \tag{12b}$$

where the subscripts  ${\bf t}$  and  ${\boldsymbol \lambda}$  represents the vector formed by stacking the transmission and the arrival rate of all users in the system. Since the transmission and the arrival rates are bounded, the second order moments in the first term of (12b) can be bounded by a constant B without affecting the optimal solution of the problem [15]. The second term in (12b) follows from the Poisson arrival process.

The expression in (12) looks similar to the WSRM formu-

lation if the weights in the WSRM problem are replaced by the numbers of backlogged packets of the corresponding users. The above approach was extended for the wireless networks in [20], in which the queues were used as weights in the WSRM formulation to determine the transmit precoders. Since the expectation is minimized by minimizing the function inside, the Q-WSRM formulation is given by

$$\max_{\mathbf{m}_{l,k,n}, \mathbf{w}_{l,k,n}} \sum_{k \in \mathcal{U}} Q_k \left( \sum_{n=1}^N \sum_{l=1}^L \log_2(1 + \gamma_{l,k,n}) \right) \tag{13a}$$
subject to.
$$\sum_{n=1}^N \sum_{k \in \mathcal{U}} \sum_{l=1}^L \operatorname{tr}\left(\mathbf{m}_{l,k,n} \mathbf{m}_{l,k,n}^H\right) \leq P_{\max} \forall b. \tag{13b}$$

To avoid excessive allocation of the resources, we include an additional rate constraint  $t_k \leq Q_k$  to address  $[x]^+$  operation in (3). The rate constrained version of the Q-WSRM, denoted by Q-WSRM extended (Q-WSRME) problem for a cellular system, is given by with the additional constraints as

$$\sum_{n=1}^{N} \sum_{l=1}^{L} \log_2(1 + \gamma_{l,k,n}) \le Q_k \, \forall k \in \mathcal{U}$$
 (14)

where the precoders are associated with  $\gamma_{l,k,n}$  defined in (2). By using the number of queued packets as the weights, the resources can be allocated to the user with more backlogged packets; this essentially results in greedy allocation.

As a special case of the problem defined in (13), we can formulate the sum rate maximization problem by setting the weights in (13a) as unity, leading to the problem as in (13) with  $Q_k = 1, \forall k \in \mathcal{U}$ . This approach provides a greedy queue minimizing allocation as compared to Q-WSRME, since the resource allocation is driven by the channel conditions in comparison to the number of queued packets as in Q-WSRME. Note that in both formulations, the resources allocated to the users are limited by the number of backlogged packets with an explicit maximum rate constraint defined by (14).

#### B. JSFRA Scheme via SINR Relaxation

The problem defined in (13) ignores the second order term arising from the Lyapunov drift minimization objective by limiting it to a constant value. In fact, using  $\ell_{q=2}$  in (5), we obtain the following objective

$$\underset{t_k}{\text{minimize}} \sum_k v_k^2 = \underset{t_k}{\text{minimize}} \sum_k Q_k^2 - 2Q_k t_k + t_k^2 \quad (15)$$

which is similar to the objective in (13). It is achieved either by removing  $t_k^2$  from (15) or when the total number of queued packets is large for all users such that  $t_k^2$  has no impact on the objective function.

By limiting  $t_k^2$  with a constant value, the Q-WSRM formulation requires an explicit rate constraint (14) to avoid overallocation of the available resources. In the proposed queue deviation formulation, the explicit rate constraint is not needed, since it is handled by the objective function (5) as discussed earlier. It makes the problem simpler and allows us to employ efficient algorithms to distribute the precoder design problem across each BS independently with minimal information exchange [11]. In contrast to the WSRM formulation, the JSFRA

and the Q-WSRME problems handle the sub-channels jointly to obtain an efficient allocation by identifying the optimal space-frequency resources for the contending users.

We present iterative algorithms to solve (6) by using alternating optimization technique in conjunction with the SCA presented [22]. The problem is to determine the transmit precoders  $\mathbf{m}_{l,k,n}$  and the receive beamformers  $\mathbf{w}_{l,k,n}$  to minimize the total number of backlogged packets in the system. The SINR expression in (2) cannot be used to formulate the problem directly due to the equality constraint. However, by using additional variables, we can relax the SINR expression in (2) by inequality constraints to solve the problem (6) as

subject to

$$\gamma_{l,k,n} \le \frac{|\mathbf{w}_{l,k,n}^{\mathbf{H}} \mathbf{H}_{l,k,n} \mathbf{m}_{l,k,n}|^2}{\beta_{l,k,n}}$$
(16b)

$$\beta_{l,k,n} \ge \widetilde{N}_0 + \sum_{(j,i) \ne (l,k)} |\mathbf{w}_{l,k,n}^{\mathbf{H}} \mathbf{H}_{b_i,k,n} \mathbf{m}_{j,i,n}|^2 \quad (16c)$$

$$\sum_{n=1}^{N} \sum_{k \in \mathcal{U}_b} \sum_{l=1}^{L} \operatorname{tr}\left(\mathbf{m}_{l,k,n} \mathbf{m}_{l,k,n}^{H}\right) \leq P_{\max} \ \forall b. (16d)$$

The SINR expression in (2) is relaxed by the inequalities (16b) and (16c). Note that (16b) is an under-estimator for SINR  $\gamma_{l,k,n}$ , and (16c) provides an upper bound for the total interference seen by user  $k \in \mathcal{U}_b$ , denoted by variable  $\beta_{l,k,n}$ . Therefore, the problem formulation in (16) is an equivalent approximation for the problem presented in (6). Note that the JSFRA formulation in (16) can be reformulated as a WSRM problem, which is known to be NP-hard [23], and therefore it belongs to the class of NP-hard problems.

In order to find a tractable solution for (16), we note that (16d) is the only convex constraint with the involved variables. Thus, we only need to deal with (16b) and (16c). We resort to the AO technique by fixing the linear receivers, to solve for the transmit beamformers. For a fixed receivers  $\mathbf{w}_{l,k,n}$ , the problem now is to find the optimal transmit beamformers  $\mathbf{m}_{l,k,n}$  which is still a challenging task. We note that for a fixed  $\mathbf{w}_{l,k,n}$ , (16c) can be written as a second-order cone (SOC) constraint. Thus, the difficulty is due to the non-convexity of the difference of convex (DC) constraint in (16b). Let us define a function,

$$f(\mathbf{u}_{l,k,n}) \triangleq \frac{|\mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{l,k,n} \mathbf{m}_{l,k,n}|^{2}}{\beta_{l,k,n}}$$
(17)

where  $\mathbf{u}_{l,k,n} \triangleq \{\mathbf{w}_{l,k,n}, \mathbf{m}_{l,k,n}, \beta_{l,k,n}\}$  is the vector which needs to be identified for the optimal allocation. Note that the function  $f(\mathbf{u}_{l,k,n})$  is convex for a fixed  $\mathbf{w}_{l,k,n}$ , since it is in fact the ratio between a quadratic form of  $\mathbf{m}_{l,k,n}$  over an affine function of  $\beta_{l,k,n}$  [24]. The nonconvex set defined by the DC constraint (16b) can be decomposed as a series of convex subsets by linearizing the convex function  $f(\mathbf{u}_{l,k,n})$  with its first order Taylor approximation around a fixed operating point  $\tilde{\mathbf{u}}_{l,k,n}$  [25], [26], also referred as SCA in [22]. By using the reduced convex subset for (16b), the problem defined in (16) can be solved at each operating point iteratively.

For this purpose, let the real and imaginary component of

the complex number  $\mathbf{w}_{l,k,n}^{\mathrm{H}}\mathbf{H}_{b_k,k,n}\mathbf{m}_{l,k,n}$  be represented by

$$p_{l,k,n} \triangleq \Re\left\{\mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b_k,k,n} \mathbf{m}_{l,k,n}\right\}$$
(18a)

$$q_{l,k,n} \triangleq \Im \left\{ \mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b_k,k,n} \mathbf{m}_{l,k,n} \right\}$$
 (18b)

and hence  $f(\mathbf{u}_{l,k,n}) = (p_{l,k,n}^2 + q_{l,k,n}^2)/\beta_{l,k,n}$ . Suppose that the current value of  $p_{l,k,n}$  and  $q_{l,k,n}$  at a specific iteration are  $\tilde{p}_{l,k,n}$  and  $\tilde{q}_{l,k,n}$ , respectively. Using first order Taylor approximation around the local point  $[\tilde{p}_{l,k,n}, \tilde{q}_{l,k,n}, \tilde{\beta}_{l,k,n}]^T$ , we can approximate (16b) by the following linear inequality

$$2\frac{\tilde{p}_{l,k,n}}{\tilde{\beta}_{l,k,n}}(p_{l,k,n} - \tilde{p}_{l,k,n}) + 2\frac{\tilde{q}_{l,k,n}}{\tilde{\beta}_{l,k,n}}(q_{l,k,n} - \tilde{q}_{l,k,n}) + \frac{\tilde{p}_{l,k,n}^2 + \tilde{q}_{l,k,n}^2}{\tilde{\beta}_{l,k,n}}\left(1 - \frac{\beta_{l,k,n} - \tilde{\beta}_{l,k,n}}{\tilde{\beta}_{l,k,n}}\right) \ge \gamma_{l,k,n}. \quad (19)$$

In summary, for fixed linear receiver  $\mathbf{w}_{l,k,n}$  and operating point  $[\tilde{p}_{l,k,n}, \tilde{q}_{l,k,n}, \tilde{\beta}_{l,k,n}]^T$ , the relaxed convex subproblem to find transmit beamformers is given by

$$\underset{\mathbf{m}_{l,k,n}, \beta_{l,k,n}}{\text{minimize}} \quad \|\tilde{\mathbf{v}}\|_{q} \tag{20a}$$

subject to  $\beta_{l,k,n} \ge \widetilde{N}_0 + \sum_{(j,i)\ne(l,k)} |\mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b_i,k,n} \mathbf{m}_{j,i,n}|^2$ (20b)

$$\begin{split} \sum_{n=1}^{N} \sum_{k \in \mathcal{U}_b} \sum_{l=1}^{L} \operatorname{tr}\left(\mathbf{m}_{l,k,n} \mathbf{m}_{l,k,n}^{\mathrm{H}}\right) &\leq P_{\max} \; \forall b \; \; \text{(20c)} \\ \text{and (19)}. \end{split}$$

Now, the optimal linear receivers for fixed transmit precoders  $\mathbf{m}_{j,i,n} \, \forall i \in \mathcal{U}, \, \forall n \in \mathcal{C}$  are obtained by minimizing (16) with respect to  $\mathbf{w}_{l,k,n}$  as

$$\begin{array}{ll}
\text{minimize} \\
\gamma_{l,k,n}, \\
\mathbf{w}_{l,k,n}, \beta_{l,k,n}
\end{array} \qquad (21a)$$

subject to 
$$\beta_{l,k,n} \ge \widetilde{N}_0 + \sum_{(j,i) \ne (l,k)} |\mathbf{w}_{l,k,n}^{\mathbf{H}} \mathbf{H}_{b_i,k,n} \mathbf{m}_{j,i,n}|^2 (21b)$$
 and (19).

Solving (21) using the KKT conditions, we obtain the following iterative expression for the receiver  $\mathbf{w}_{l,k,n}^{o}$  as

$$\mathbf{A}_{l,k,n} = \sum_{(j,i)\neq(l,k)} \mathbf{H}_{b_i,k,n} \mathbf{m}_{j,i,n} \mathbf{m}_{j,i,n}^{\mathbf{H}} \mathbf{H}_{b_i,k,n}^{\mathbf{H}} + N_0 \mathbf{I}_{N_R}$$
(22a)

$$\mathbf{w}_{l,k,n}^{(i)} = \left(\frac{\tilde{\beta}_{l,k,n} \mathbf{m}_{l,k,n}^{H} \mathbf{H}_{b_k,k,n}^{H} \mathbf{w}_{l,k,n}^{(i-1)}}{\|\mathbf{w}_{l,k,n}^{(i-1)} \mathbf{H}_{b_k,k,n} \mathbf{m}_{l,k,n}\|^2}\right) \mathbf{A}_{l,k,n}^{-1} \mathbf{H}_{b_k,k,n} \mathbf{m}_{l,k,n} (22b)$$

where  $\mathbf{w}_{l,k,n}^{(i-1)}$  is the receive beamformer from the previous iteration, upon which the linear relaxation is performed for the nonconvex DC constraint in (16b), as used in the formulation (21). The optimal receiver  $\mathbf{w}_{l,k,n}^o$  is obtained by either iterating (22b) until convergence or for a fixed number of iterations. Note that the receiver has no explicit relation with the choice of  $\ell_q$  norm used in the objective function. The dependency is implied by the transmit precoders  $\mathbf{m}_{l,k,n}$ , which depend on the value of the exponent q.

<sup>&</sup>lt;sup>4</sup>Note that  $p_{l,k,n}$  and  $q_{l,k,n}$  are just symbolic notation and not the newly introduced optimization variables. In CVX [27], for example, we declare  $p_{l,k,n}$  and  $q_{l,k,n}$  with the 'expression' qualifier.

It can be seen that the optimal receiver in (22b) is in fact a scaled version of the MMSE receiver, which is given by

$$\mathbf{R}_{l,k,n} = \sum_{i \in \mathcal{U}} \sum_{j=1}^{L} \mathbf{H}_{b_i,k,n} \mathbf{m}_{j,i,n} \mathbf{m}_{j,i,n}^{\mathbf{H}} \mathbf{H}_{b_i,k,n}^{\mathbf{H}} + N_0 \mathbf{I}_{N_R}$$
(23a)

$$\mathbf{w}_{l,k,n} = \mathbf{R}_{l,k,n}^{-1} \,\mathbf{H}_{b_k,k,n} \,\mathbf{m}_{l,k,n}. \tag{23b}$$

Since the scaling present in the optimal receiver (22b) has no impact on the received SINRs, the MMSE receiver in (23b) can also be used without compromising the performance or the convergence behavior.

#### Algorithm 1: Algorithm of JSFRA scheme

```
Input: a_k, Q_k, \mathbf{H}_{b,k,n}, \forall b \in \mathcal{B}, \forall k \in \mathcal{U}, \forall n \in \mathcal{N}
Output: \mathbf{m}_{l,k,n} and \mathbf{w}_{l,k,n} \forall l \in \{1,2,\ldots,L\}
Initialize: i = 0 and transmit precoders \tilde{\mathbf{m}}_{l,k,n} randomly
               satisfying the total power constraint (6b)
update \mathbf{w}_{l,k,n}, \mathbf{u}_{l,k,n} using (23b) and (19) with \mathbf{m}_{l,k,n}
     initialize j = 0
     repeat
          solve for the transmit precoders \mathbf{m}_{l,k,n} using (20)
          update the constraint set (19) with \mathbf{u}_{l,k,n} and
           \mathbf{m}_{l,k,n} using (18)
          j = j + 1
     until SCA convergence or j \geq J_{\max}
     update the receive beamformers \mathbf{w}_{l,k,n} using (21) or
     (23b) with the updated precoders \mathbf{m}_{l,k,n}
     i = i + 1
until Queue convergence or i \geq I_{\text{max}}
```

The proposed subproblems in (20) and (21) are solved in an iterative manner by updating the operating point from the previous iteration. The iterative algorithm is referred to as queue minimizing JSFRA scheme with a per BS power constraint, and it is outlined in Algorithm 1. The iterative procedure repeats until the improvement on the objective is less than a predetermined tolerance parameter or the maximum number of iterations is reached. Instead of initializing  $\mathbf{u}_{l,k,n}$ arbitrarily to a feasible point, transmit precoders can also be initialized with some feasible point  $\tilde{\mathbf{m}}_{l,k,n}$ , which is then used to find  $\mathbf{u}_{l,k,n}$  as seen in Algorithm 1. For a fixed receive beamformer  $\mathbf{w}_{l,k,n}$ , the SCA iteration is carried out until convergence or for predefined number of iterations, say,  $J_{\text{max}}$ for the optimal transmit precoders  $\mathbf{m}_{l,k,n}$ . Next, the receive beamformers are updated based on either (22b) or (23b) using the fixed transmit precoders  $\mathbf{m}_{l,k,n}$ . This procedure is carried out until convergence of the queue deviation or a fixed number of iterations  $I_{\text{max}}$  as outlined in Algorithm 1. The convergence proof is discussed in Appendix A.

#### C. JSFRA Scheme via MSE Reformulation

In the second method, we solve the JSFRA problem by exploiting the relation between the MSE and the achievable SINR when the MMSE receivers are used at the user terminals [4], [5]. The MSE  $\epsilon_{l,k,n}$ , for a data symbol  $d_{l,k,n}$  is given by

$$\mathbb{E}\left[\left(d_{l,k,n} - \hat{d}_{l,k,n}\right)^{2}\right] = \left|1 - \mathbf{w}_{l,k,n}^{H} \mathbf{H}_{b_{k},k,n} \mathbf{m}_{l,k,n}\right|^{2} + \sum_{(j,i)\neq(l,k)} \left|\mathbf{w}_{l,k,n}^{H} \mathbf{H}_{b_{i},k,n} \mathbf{m}_{j,i,n}\right|^{2} + \widetilde{N}_{0} = \epsilon_{l,k,n} \quad (24)$$

where  $\hat{d}_{l,k,n}$  is the estimate of the transmitted symbol. Using the MMSE receive beamformer (23b) in the MSE expression (24) and in the SINR expression (2), we can arrive at the following relation between the MSE and the SINR as

$$\epsilon_{l,k,n} = (1 + \gamma_{l,k,n})^{-1}.$$
 (25)

The above equivalence is valid only if the receivers are based on the MMSE criterion. Using the equivalence in (25), the WSRM objective can be reformulated as the weighted minimum mean squared error (WMMSE) equivalent to obtain the precoders for the MU-MIMO scenario as discussed in [5]–[7]. Note that the receiver is invariably based on the MMSE criterion irrespective of the  $\ell_q$  norm used in the objective function to obtain the optimal transmit precoders  $\mathbf{m}_{l,k,n}$ .

Let  $v_k' = Q_k - \sum_{n=1}^N \sum_{l=1}^L t_{l,k,n}$  denote the queue deviation corresponding to user k and  $\tilde{v}_k' \triangleq a_k^{1/q} v_k'$  represents the weighted equivalent. By using the relaxed MSE expression in (24), the problem in (6) can be expressed as

$$\underset{\substack{t_{l,k_n}, \mathbf{m}, \mathbf{w}, t_{k,n}, \\ t_{l+1}, \dots, \mathbf{w}, \mathbf{w}, t_{l+n}, \\ t_{l+1}, \dots, \mathbf{w}, t_{l+n}}}{\text{minimize}}, \quad \|\tilde{\mathbf{v}}'\|_q \tag{26a}$$

subject to 
$$t_{l,k,n} \leq -\log_2(\epsilon_{l,k,n})$$
 (26b)
$$\sum_{(j,i)\neq(l,k)} \left| \mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b_i,k,n} \mathbf{m}_{j,i,n} \right|^2 + \widetilde{N}_0$$

$$+\left|1-\mathbf{w}_{l,k,n}^{\mathrm{H}}\mathbf{H}_{b_{k},k,n}\mathbf{m}_{l,k,n}\right|^{2} \le \epsilon_{l,k,n}$$
 (26c)

$$\sum_{n=1}^{N} \sum_{k \in \mathcal{U}_b} \sum_{l=1}^{L} \operatorname{tr}\left(\mathbf{m}_{l,k,n} \mathbf{m}_{l,k,n}^{H}\right) \leq P_{\max} \ \forall b. (26d)$$

The alternative MSE formulation given by (26) is non-convex even for the fixed  $\mathbf{w}_{l,k,n}$  due to the constraint (26b), which is in fact a DC constraint. We resort to the SCA approach [22] by relaxing the constraint by a sequence of convex subsets using first order Taylor series approximation around a fixed MSE point  $\tilde{\epsilon}_{l,k,n}$  as

$$-\log_2(\tilde{\epsilon}_{l,k,n}) - \frac{(\epsilon_{l,k,n} - \tilde{\epsilon}_{l,k,n})}{\log(2)\,\tilde{\epsilon}_{l,k,n}} \ge t_{l,k,n} \tag{27}$$

Using the above approximation for the rate constraint, the problem defined in (26) is solved for optimal transmit precoders  $\mathbf{m}_{l,k,n}$ , MSEs  $\epsilon_{l,k,n}$ , and the user rates over each sub-channel  $t_{l,n,k}$  for a fixed receive beamformers. The optimization subproblem to find the transmit precoders for a fixed receive beamformers  $\mathbf{w}_{l,k,n}$  is given by

$$\begin{array}{ll}
\underset{t_{l,k,n},\mathbf{m}_{l,k,n},\epsilon_{l,k,n}}{\text{minimize}} & \|\tilde{\mathbf{v}}'\|_q \\
\text{subject to} & (26c), (26d), \text{ and } (27).
\end{array} (28a)$$

The optimal transmit precoders for a fixed receivers are obtained by solving the subproblem (28) iteratively by updating the fixed MSE point  $\tilde{\epsilon}_{l,k,n}$  with  $\epsilon_{l,k,n}$  from the previous iteration until termination as discussed in Section III-B. The convergence analysis follows the discussions in Appendix A.

#### D. Reduced Complexity Spatial Resource Allocation (SRA)

The complexity of the JSFRA algorithm scales quickly with the number of sub-channels, since the complexity of the interior point method, which is used to solve the problem, increases quickly with the problem size. Thus, we can use the decomposition methods presented in [10], [11] to overcome this complexity by designing precoders for each sub-channels independently with minimal information exchange.

As an alternative sub-optimal solution, we present queue minimizing spatial resource allocation (SRA), which solves for the precoders using JSFRA formulation for a specific subchannel i with a fixed transmit power  $P_{\max,i}$ . The sharing of power can be equal or based on a predetermined pattern as in partial frequency reuse for the sub-channels as given by

$$\sum_{i=1}^{N} P_{\max,i} = P_{\max}.$$
 (29)

Even though N sub-channels are present at any given scheduling instant, precoders are computed for each sub-channel in a sequential manner with the sub-channel specific total power constraint  $P_{\max,i}$  and the number of backlogged packets. Let  $Q_{k,i}$  be the number of backlogged packets associated with user k before solving for the precoders specific to the sub-channel i. Since the precoder design is sequential, i.e, the precoders are designed for sub-channels [0,i-1] before  $i^{\text{th}}$  sub-channel, the number of backlogged packets for the initial sub-channel is initialized as  $Q_{k,1} = Q_k$ . The queues associated with the consecutive sub-channels are given by

$$Q_{k,i+1} = \max\left(Q_k - \sum_{j=1}^{i} \sum_{l=1}^{L} t_{l,k,j}, 0\right) \ \forall \ k \in \mathcal{U}$$
 (30)

where  $t_{l,k,j}$  denotes the rate corresponding to the user k on the  $j^{\rm th}$  sub-channel and  $l^{\rm th}$  spatial stream. Note that the proposed scheme is sensitive to the order of the sub-channel selection due to the sequential precoder design for each sub-channel. However, the SRA approach provides faster convergence in contrast to the JSFRA formulation due to the substantial reduction in the optimization variables for each sub-channel problem. As the number of users in the system increases, the SRA formulation will be insensitive to the sub-channel ordering due to the multi-user diversity.

#### IV. DISTRIBUTED SOLUTIONS

The distributed precoder designs for the proposed JSFRA scheme are discussed in this section. The convex formulation in (20) or (28) requires a centralized controller to perform the precoder design for all users belonging to the coordinating BSs. In order to design the precoders independently at each BS with the minimal information exchange via backhaul, iterative decentralization methods are addressed. In particular, the primal decomposition and the ADMM based dual decomposition approaches are considered.

Let us consider the convex subproblem with the fixed receive beamformers  $\mathbf{w}_{l,k,n}$  presented in (20) based on the Taylor series approximation for the nonconvex constraint. The following discussions are equally valid for the MSE based

solution outlined in (28) as well. Since the objective of (20) can be decoupled across each BS, the centralized problem can be equivalently written as

$$\underset{\gamma_{l,k,n},\mathbf{m}_{l,k,n},\beta_{l,k,n}}{\text{minimize}} \qquad \sum_{b \in \mathcal{B}} \|\tilde{\mathbf{v}}_b\|_q \tag{31a}$$

subject to 
$$(20b) - (20d)$$
 (31b)

where  $\tilde{\mathbf{v}}_b$  denotes the vector of weighted queue deviation corresponding to users  $k \in \mathcal{U}_b$ .

To begin with, let  $\bar{\mathcal{B}}_b$  denote the set  $\mathcal{B}\setminus\{b\}$  and  $\bar{\mathcal{U}}_b$  represents the set  $\mathcal{U}\setminus\mathcal{U}_b$ . Following an approach similar to the one presented in [12], [13], the coupling constraint (20b) or (26c) can be expressed by grouping the interference from each BS in  $\bar{\mathcal{B}}_{bb}$  as

$$\widetilde{N}_{0} + \sum_{j=1, j \neq l}^{L} |\mathbf{w}_{l,k,n}^{H} \mathbf{H}_{b_{k},k,n} \mathbf{m}_{j,k,n}|^{2} + \sum_{b \in \overline{\mathcal{B}}_{b_{k}}} \zeta_{l,k,n,b}$$

$$+ \sum_{i \in \mathcal{U}_{b_{k}} \setminus \{k\}} \sum_{j=1}^{L} |\mathbf{w}_{l,k,n}^{H} \mathbf{H}_{b_{k},k,n} \mathbf{m}_{j,i,n}|^{2} \leq \beta_{l,k,n} \quad (32)$$

where  $\zeta_{l,k,n,b}$  is the total interference caused by the transmission of BS b to user  $k \in \mathcal{U}_{b_k}$  in the spatial stream l and sub-channel n. It is given by the following upper bound as

$$\zeta_{l,k,n,b} \ge \sum_{i \in \mathcal{U}_b} \sum_{j=1}^{L} |\mathbf{w}_{l,k,n}^{\mathbf{H}} \mathbf{H}_{b,k,n} \mathbf{m}_{j,i,n}|^2 \, \forall b \in \bar{\mathcal{B}}_{b_k}. \tag{33}$$

The decentralization is achieved by decomposing the original convex problem in (31) to a parallel iterative subproblems coordinated by either primal or dual decomposition update. The coupling variables are updated in each iteration by exchanging limited information among the subproblems. Before proceeding further, let  $\bar{\zeta}_b$  be the vector formed by stacking interference terms (33) from the neighboring BSs to the users of BS b and  $\hat{\zeta}_b$  be the stacked interference terms caused by BS b to all users in the neighboring BSs  $\bar{\mathcal{B}}_b$ , represented as

$$\bar{\zeta}_{b} = \left[\zeta_{l,k,n,\bar{\mathcal{B}}_{b}(1)}, \dots, \zeta_{l,k,n,\bar{\mathcal{B}}_{b}(|\bar{\mathcal{B}}_{b}|)}\right]^{\mathrm{T}}, \forall k \in \mathcal{U}_{b}$$
(34a)
$$\hat{\zeta}_{b} = \left[\zeta_{l,\bar{\mathcal{U}}_{b}(1),n,b}, \zeta_{l,\bar{\mathcal{U}}_{b}(2),n,b}, \dots, \zeta_{l,\bar{\mathcal{U}}_{b}(|\bar{\mathcal{U}}_{b}|),n,b}\right]^{\mathrm{T}}$$
(34b)

Let us define the vector  $\zeta_b$ , formed by stacking the interference terms corresponding to the BS b as

$$\boldsymbol{\zeta}_b = \left[\hat{\boldsymbol{\zeta}}_b^{\mathrm{T}}, \bar{\boldsymbol{\zeta}}_b^{\mathrm{T}}\right]^{\mathrm{T}}.$$
 (35)

Since the decentralization solution is an iterative procedure, we represent the  $i^{\rm th}$  iteration index as  $x^{(i)}$ . Let  $\zeta_b(b_k)$  denote the interference terms corresponding to BS  $b_k$  in BS b as

$$\boldsymbol{\zeta}_b(b_k) = \left[\zeta_{l,\mathcal{U}_b(1),n,b_k}, \dots, \zeta_{l,\mathcal{U}_b(|\mathcal{U}_b|),n,b_k}\right]. \tag{36}$$

To decentralize the problem in (31), the BS specific vector  $\zeta_b$  in (35), which are relevant for the BS b, can either be fixed or treated as a variable in accordance to the decomposition method. To decouple the precoder design across BSs, the equivalent downlink channel  $\mathbf{w}_{l,k,n}^{\mathrm{H}}\mathbf{H}_{b,k,n}, \forall k \in \mathcal{U}$  are to be known at each BS b through the precoded uplink pilots from all the users in the system, where the precoders are the MMSE

receiver  $\mathbf{w}_{l,k,n}$  evaluated at the user. Similarly, to update the MMSE receivers at each user k, the equivalent channels  $\mathbf{H}_{b,k,n}\mathbf{m}_{l,k',n}, \forall k' \in \mathcal{U}_b, \forall b \in \mathcal{B}$  need to be known through the user specific precoded downlink pilots precoded with the updated transmit beamformers  $\mathbf{m}_{l,k,n}, \forall k \in \mathcal{B}$  evaluated at the BS b by using the equivalent downlink channel that includes the updated MMSE receivers of all the users, as in [28].

#### A. Primal Decomposition

In the primal decomposition, the convex problem in (31) is solved for the optimal transmit precoders in an iterative manner for a fixed BS specific interference term  $\zeta_{b_k}$  using the master-slave model [12]. The slave subproblem is solved in each BS for the optimal transmit precoders only for the associated users by assuming fixed interference terms  $\zeta_{b_k}^{(i)}$  in each  $i^{\text{th}}$  iteration. Upon finding the optimal associated transmit precoders by each slave subproblems, the master problem is used to update the BS specific interference terms  $\zeta_{b_k}^{(i+1)}$ for the next iteration by using dual variables corresponding to the interference constraint (32) as discussed in [12]. In this manner, the interference variables are updated until the global consensus is obtained. The primal approach is similar to the minimum power precoder design presented in [12]. Note that the master problem treats  $\zeta_b$  as a variable and the slave subproblems assumes it to be a constant for each iteration to find the transmit precoders.

#### B. Alternating Directions Method of Multipliers (ADMM)

The ADMM approach is used to decouple the precoder design across multiple BSs to solve the convex problem in (31). The ADMM is preferred over the dual decomposition (DD) approach in [13] for its robustness and improved convergence behavior [11]. In contrast to the primal decomposition, the ADMM approach relaxes the interference constraints by including in the objective function of each subproblem with a penalty pricing [10], [11]. Similar approach for the precoder design in the minimum power context is considered in [29].

Using the formulation presented in [11], [29], we can write the BS b specific ADMM subproblem for the  $i^{\rm th}$  iteration as

$$\underset{\substack{\gamma_{l,k,n},\mathbf{m}_{l,k,n},\zeta_{b}}{\beta_{l,k,n},\zeta_{b}}}{\text{minimize}} \|\tilde{\mathbf{v}}_{b}\|_{q} + \boldsymbol{\nu}_{b}^{(i) \mathrm{T}} \left(\boldsymbol{\zeta}_{b} - \boldsymbol{\zeta}_{b}^{(i)}\right) + \frac{\rho}{2} \|\boldsymbol{\zeta}_{b} - \boldsymbol{\zeta}_{b}^{(i)}\|^{2} \tag{37a}$$

subject to 
$$\sum_{n=1}^{N} \sum_{k \in \mathcal{U}_h} \sum_{l=1}^{L} \operatorname{tr}\left(\mathbf{m}_{l,k,n} \mathbf{m}_{l,k,n}^{H}\right) \leq P_{\max}$$
 (37b)

$$\sum_{\bar{b}\in\bar{\mathcal{B}}_{b}} \zeta_{l,k,n,\bar{b}} + \sum_{\{\bar{l},\bar{k}\}\neq\{l,k\}} |\mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b,k,n} \mathbf{m}_{\bar{l},\bar{k},n}|^{2} + \widetilde{N}_{0} \leq \beta_{l,k,n}$$

$$\sum_{k \in \mathcal{U}_b} \sum_{l=1}^{L} |\mathbf{w}_{\bar{l},\bar{k},n}^{\mathrm{H}} \mathbf{H}_{b,\bar{k},n} \mathbf{m}_{l,k,n}|^2 \le \zeta_{\bar{l},\bar{k},n,b} \ \forall \bar{k} \in \bar{\mathcal{U}}_b \ \forall n$$
 (37d) and (19)

where  $\zeta_b^{(i)}$  denotes the interference vector updated from the earlier iteration and  $\boldsymbol{\nu}_b^{(i)}$  represents the dual vector corresponding to the equality constraint at the  $i^{\mathrm{th}}$  iteration as

$$\zeta_b = \zeta_b^{(i)}.\tag{38}$$

#### Algorithm 2: Distributed JSFRA scheme using ADMM

**Input**:  $a_k$ ,  $Q_k$ ,  $\mathbf{H}_{b,k,n}$ ,  $\forall b \in \mathcal{B}$ ,  $\forall k \in \mathcal{U}$ ,  $\forall n \in \mathcal{N}$ 

```
Output: \mathbf{m}_{l,k,n} and \mathbf{w}_{l,k,n} \forall l
Initialize: i = 0 and \mathbf{m}_{l,k,n} randomly satisfying (37b)
update \mathbf{w}_{l,k,n} with (23b) and \tilde{\mathbf{u}}_{l,k,n} using (16c) and (18)
initialize the interference vectors \boldsymbol{\zeta}_b^{(0)} = \mathbf{0}^{\mathrm{T}}, \forall b \in \mathcal{B}
initialize the dual vectors \boldsymbol{\nu}_b^{(0)} = \boldsymbol{0}^{\mathrm{T}}, \forall b \in \mathcal{B}
foreach BS b \in \mathcal{B} do
      repeat
           initialize j = 0
            repeat
                  solve for \mathbf{m}_{l,k,n} and \boldsymbol{\zeta}_b with (37) using \boldsymbol{\zeta}_b^{(j)}
                 exchange \zeta_b among BSs in \mathcal{B} update \zeta_b^{(j+1)} and \nu_b^{(j+1)} using (39) and (40) j=j+1
            until convergence or j \geq J_{\max}
            precoded downlink pilot transmission with \mathbf{m}_{l,k,n}
            update \mathbf{w}_{l,k,n} to all BSs in \mathcal{B} using precoded
            uplink pilots [28]
            update \tilde{\mathbf{u}}_{l,k,n} using (16c) and (18) for SCA point
            or \tilde{\epsilon}_{l,k,n} using (26c) for MSE operating point
      until convergence or i \geq I_{\max}
end
```

Upon solving (37) for  $\zeta_b \forall b$  in the  $i^{\rm th}$  iteration, the next iterate is updated by exchanging the corresponding interference terms between two BSs b and  $b_k$  as

$$\zeta_{b_k}(b)^{(i+1)} = \zeta_b(b_k)^{(i+1)} = \frac{\zeta_b(b_k) + \zeta_{b_k}(b)}{2}.$$
(39)

The dual vector for the next iteration is updated by using the subgradient search to maximize the dual objective as

$$\nu_b^{(i+1)} = \nu_b^{(i)} + \rho \left( \zeta_b - \zeta_b^{(i+1)} \right) \tag{40}$$

where step size parameter *ρ* is chosen in accordance with [11] to depend on the system model under consideration. The convergence rate of the distributed algorithm is susceptible to the choice of step size parameter *ρ*. For our simulation models, we consider step size parameter *ρ* = 2. The above iteration is performed until convergence or for certain accuracy in the variation of the objective value between two consecutive updates. The distributed precoder design using the ADMM approach is shown in Algorithm 2. The convergence analysis (37c) of the distributed algorithms are discussed in Appendix B.

#### C. Decomposition via KKT Conditions for MSE Formulation

In this section, we discuss an alternative way to decentralize the precoder design across the coordinating BSs in  $\mathcal B$  based on the MSE reformulation method discussed in Section III-C. In contrast to Section IV-A and IV-B, the problem is solved using the KKT conditions in which the transmit precoders, receive beamformers and the subgradient updates are performed at the same instant to minimize the global queue deviation objective with few number of iterations. The proposed methods

in this section provide algorithms that can be of practical importance owing to the limited signaling requirements. We consider an idealized TDD system due to the knowledge of complete channel information at the transmitter. Similar work has been considered for the WSRM problem with minimum rate constraints in [8], [9]. Since the formulations in [8], [9] are similar to that of the Q-WSRME scheme with an additional maximum rate constraint (14), it requires explicit dual variables to handle the maximum rate constraint, thereby making the problem difficult to solve in an iterative manner.

In the proposed JSFRA formulation, the maximum rate constraints are implicitly handled by the objective function without the need of explicit constraints. However, the KKT conditions cannot be formulated due to the non-differential objective function. The non-differentiability is due to the absolute value operator present in the norm function. In order to make the objective function differentiable, we consider the following two cases for which the absolute operator can be ignored without affecting the optimal solution, namely,

- when the exponent q is even, or
- when the number of backlogged packets of each user is large enough, i.e,  $Q_k \gg \sum_{n=1}^N \sum_{l=1}^L t_{l,k,n}$  to ignore the absolute operator and queues in the first place as well.

With the assumption of either one of the above conditions to be true, the problem in (28) can be written as

$$\underset{\substack{t_{l,k,n}, \mathbf{m}_{l,k,n}, \\ e_{l,k,n}, \mathbf{w}_{l,k,n}, \\ \text{subject to}}}{\text{minimize}} \sum_{b \in \mathcal{B}} \sum_{k \in \mathcal{U}_b} a_k \left( Q_k - \sum_{n=1}^N \sum_{l=1}^L t_{l,k,n} \right)^q \tag{41a}$$

$$\alpha_{l,k,n} : \left| 1 - \mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b_k,k,n} \mathbf{m}_{l,k,n} \right|^2 + \widetilde{N}_0$$

$$+ \sum_{(x,y) \neq (l,k)} \left| \mathbf{w}_{l,k,n}^{\mathrm{H}} \mathbf{H}_{b_y,k,n} \mathbf{m}_{x,y,n} \right|^2 \leq \epsilon_{l,k,n}$$
(41b)

$$\sigma_{l,k,n} : \log_2(\tilde{\epsilon}_{l,k,n}) + \frac{(\epsilon_{l,k,n} - \tilde{\epsilon}_{l,k,n})}{\log(2)\tilde{\epsilon}_{l,k,n}} \le -t_{l,k,n} \quad (41c)$$

$$\delta_b: \sum_{n=1}^{N} \sum_{k \in \mathcal{U}_b} \sum_{l=1}^{L} \operatorname{tr}\left(\mathbf{m}_{l,k,n} \mathbf{m}_{l,k,n}^{H}\right) \le P_{\max} \ \forall b \ \ (41d)$$

where  $\alpha_{l,k,n}$ ,  $\sigma_{l,k,n}$  and  $\delta_b$  are the dual variables corresponding to the constraints defined in (41b), (41c) and (41d).

The problem in (41) is solved using the KKT expressions, obtained by the derivative of the Lagrangian function w.r.t the primal and the dual variables, complementary slackness, and the primal, dual feasibility requirements as shown in Appendix C. Upon solving, we obtain the iterative solution as

$$\mathbf{m}_{l,k,n}^{(i)} = \left(\sum_{x \in \mathcal{U}} \sum_{y=1}^{L} \alpha_{y,x,n}^{(i-1)} \mathbf{H}_{b_{k},x,n}^{H} \mathbf{w}_{y,x,n}^{(i-1)} \mathbf{w}_{y,x,n}^{H (i-1)} \mathbf{H}_{b_{k},x,n} \right. \\ + \left. \delta_{b} \mathbf{I}_{N_{T}} \right)^{-1} \alpha_{l,k,n}^{(i-1)} \mathbf{H}_{b_{k},k,n}^{H} \mathbf{w}_{l,k,n}^{(i-1)} \qquad (42a)$$

$$\mathbf{w}_{l,k,n}^{(i)} = \left(\sum_{x \in \mathcal{U}} \sum_{y=1}^{L} \mathbf{H}_{b_{x},k,n} \mathbf{m}_{y,x,n}^{(i)} \mathbf{m}_{y,x,n}^{H (i)} \mathbf{H}_{b_{x},k,n}^{H} \right. \\ + \left. \mathbf{I}_{N_{R}} \right)^{-1} \mathbf{H}_{b_{k},k,n} \mathbf{m}_{l,k,n}^{(i)}$$

$$\epsilon_{l,k,n}^{(i)} = \left. \left| 1 - \mathbf{w}_{l,k,n}^{H (i)} \mathbf{H}_{b_{k},k,n} \mathbf{m}_{l,k,n}^{(i)} \right|^{2} + N_{0} \left\| \mathbf{w}_{l,k,n}^{(i)} \right\|^{2} \right.$$

$$+\sum_{(x,y)\neq(l,k)} \left| \mathbf{w}_{l,k,n}^{\mathrm{H}(i)} \mathbf{H}_{b_y,k,n} \mathbf{m}_{x,y,n}^{(i)} \right|^2$$
(42c)

$$t_{l,k,n}^{(i)} = -\log_2(\epsilon_{l,k,n}^{(i-1)}) - \frac{\left(\epsilon_{l,k,n}^{(i)} - \epsilon_{l,k,n}^{(i-1)}\right)}{\log(2)\epsilon_{l,k,n}^{(i-1)}}$$
(42d)

$$\sigma_{l,k,n}^{(i)} = \left[ \frac{a_k q}{\log(2)} \left( Q_k - \sum_{n=1}^{N} \sum_{l=1}^{L} t_{l,k,n}^{(i)} \right)^{(q-1)} \right]^+ \tag{42e}$$

$$\alpha_{l,k,n}^{(i)} = \alpha_{l,k,n}^{(i-1)} + \rho \left( \frac{\sigma_{l,k,n}^{(i)}}{\epsilon_{l,k,n}^{(i)}} - \alpha_{l,k,n}^{(i-1)} \right)$$
(42f)

Since the dual variables  $\alpha^{(i)}$  and  $\sigma^{(i)}$  are interdependent in (42), one has to be fixed to optimize for the other. So,  $\alpha^{(i)}$  is fixed to evaluate  $\sigma^{(i)}$  using (42). At each iteration, the dual variables  $\alpha^{(i)}$  are linearly interpolated with any point between the fixed iterate  $\alpha^{(i-1)}$  and  $\frac{\sigma^{(i)}}{\epsilon^{(i)}}$  using a step size  $\rho \in (0,1)$ . The choice of  $\rho$  depends on the system model and it affects the convergence behavior. It reduces the oscillations in the objective function when  $\sigma^{(i)}$  is negative due to over-allocation.

objective function when  $\sigma^{(i)}$  is negative due to over-allocation. When the allocated rate  $t_k^{(i-1)}$  is greater than the number of queued packets  $Q_k$  for a user k, the corresponding dual variable  $\sigma^{(i)}$  will be negative and due to the projection operator  $[x]^+$  in (42e), it will be zero, thereby forcing  $\alpha_k^{(i)} < \alpha_k^{(i-1)}$  as in (42f). Once  $\alpha_k^{(i)}$  is reduced, the precoder weight in (42a) is lowered to make the rate  $t_k^{(i)} < t_k^{(i-1)}$  eventually. The choice of  $\rho$  is susceptible to the system model under consideration, which affects the convergence speed of the iterative algorithm. In all simulations, the step size  $\rho$  is fixed to 0.1 irrespectively.

The KKT expressions in (42) are solved in an iterative manner by initializing the transmit and the receive beamformers  $\mathbf{m}_{l,k,n}$ ,  $\mathbf{w}_{l,k,n}$  with the single user beamforming and the MMSE vectors. The dual variable  $\alpha$ 's are initialized with ones to have equal priorities to all the users in the system. Then the transmit and the receive beamformers are evaluated using the expressions in (42). The transmit precoder in (42a) depends on the BS specific dual variable  $\delta_b$ , which can be found by bisection search satisfying the total power constraint (41d). Note that the fixed SCA operating point is given by  $\tilde{\epsilon}_{l,k,n} = \epsilon_{l,k,n}^{(i-1)}$ , which is considered in the expression (42).

To obtain a practical distributed precoder design, we assume that each BS b knows the corresponding equivalent channel  $\mathbf{w}_{l,k,n}^{\mathrm{H}}\mathbf{H}_{b,k,n}, \forall k \in \mathcal{U}$ , which includes the receivers  $\mathbf{w}_{l,k,n}$ , through precoded uplink pilot signaling. We extend the decentralization methods discussed in [28], for the current problem as follows. After receiving the updated transmit precoders from all BSs in  $\mathcal{B}$ , each user evaluates the MMSE receiver in (42b) and notify them to the BSs via uplink precoded pilots. On receiving pilot signals, BSs update the MSE in (24) as

$$\epsilon_{l,k,n}^{(i)} = 1 - \mathbf{w}_{l,k,n}^{(i)H} \mathbf{H}_{b_k,k,n} \mathbf{m}_{l,k,n}^{(i)}.$$
 (43)

Using the current MSE value,  $t_{l,k,n}^{(i)}, \sigma_{l,k,n}^{(i)}$ , and  $\alpha_{l,k,n}^{(i)}$  are evaluated using (42d), (42e) and (42f), and the updated dual variables  $\alpha_{l,k,n}$  are exchanged between the BSs to evaluate the transmit precoders  $\mathbf{m}_{l,k,n}^{(i+1)}$  for the next iteration. The SCA operating point is also updated with the current MSE value.

To avoid the back-haul exchanges between the BSs, as an alternative approach, users perform all processing required and BSs will update the precoders based on the feedback

Algorithm 3: KKT approach for the JSFRA scheme

```
Input: a_k, Q_k, \mathbf{H}_{b,k,n}, \forall b \in \mathcal{B}, \forall k \in \mathcal{U}, \forall n \in \mathcal{N}
Output: \mathbf{m}_{l,k,n} and \mathbf{w}_{l,k,n} \forall l \in \{1,2,\ldots,L\}
Initialize: i=1, \mathbf{w}_{l,k,n}^{(0)}, \tilde{\epsilon}_{l,k,n} randomly, dual variables
                  \alpha_{l,k,n}^{(0)} = 1, and I_{\text{max}} for certain value
foreach BS b \in \mathcal{B} do
      initialize i = 0
      repeat
             update \mathbf{m}_{l,k,n}^{(i)} using (42a), and perform precoded
             downlink pilot transmission
             find \mathbf{w}_{l,k,n}^{(i)} using (42b) at each user
            evaluate \epsilon_{l,k,n}^{(i)}, t_{l,k,n}^{(i)}, \sigma_{l,k,n}^{(i)} and \alpha_{l,k,n}^{(i)} using (42c) and (42d), (42e) and (42f) at each user with
             the updated \mathbf{w}_{l,k,n}^{(i)}
             using precoded uplink pilots, \mathbf{m}_{l,k,n}^{(i)} and \alpha_{l,k,n}^{(i)}
             are notified to all BSs in \ensuremath{\mathcal{B}}
             i = i + 1
      until until convergence or i \geq I_{\max}
end
```

information from the users. Upon receiving the transmit precoders from BSs, each user will update the receive beamformer  $\mathbf{w}_{l,k,n}$ , the MSE  $\epsilon_{l,k,n}$ , and the dual variables  $\lambda_{l,k,n}$  and  $\alpha_{l,k,n}$ . The updated  $\alpha_{l,k,n}$  and  $\mathbf{w}_{l,k,n}$  are notified to the BSs using two separate precoded uplink pilot symbols with  $\tilde{\mathbf{w}}_{l,k,n}^{(i)} = \sqrt{\alpha_{l,k,n}^{(i)}}\mathbf{w}_{l,k,n}^{*(i)}$  and  $\bar{\mathbf{w}}_{l,k,n}^{(i)} = \alpha_{l,k,n}^{(i)}\mathbf{w}_{l,k,n}^{*(i)}$  as the precoders. On receiving the precoded uplink pilots, each BS use the effective channel  $\mathbf{H}_{b,k,n}^{\mathrm{T}}\tilde{\mathbf{w}}_{l,k,n}^{(i)}$  and  $\mathbf{H}_{b,k,n}^{\mathrm{T}}\bar{\mathbf{w}}_{l,k,n}^{(i)}$  in (42a) to update the transmit precoders, where  $\mathbf{x}^*$  is the complex conjugate of  $\mathbf{x}$ . Algorithm 3 outlines a practical way of updating the transmit and the receive beamformers by using over-the-air (OTA) signaling with the precoded pilots for the KKT based MSE reformulated JSFRA problem. The convergence of the algorithm is discussed in Appendix B.

#### V. SIMULATION RESULTS

The simulations carried out in this work consider the path loss (PL) varying uniformly across all users in the system with the channels drawn from the *i.i.d.* samples. The queues are generated based on the Poisson process with the average values specified in each section presented.

#### A. Centralized Solutions

We discuss the performance of the centralized algorithms in Section III for some system configurations. To begin with, we consider a single cell single-input single-output (SISO) model operating at 10 dB signal-to-noise ratio (SNR) with K=3 users sharing N=3 sub-channel resources. The number of packets waiting at the transmitter for each user is given by  $Q_k=4,8$  and 4 bits, respectively.

Table I tabulate the channels of the users over each subchannel followed by the rates assigned by three different algorithms, Q-WSRME allocation, JSFRA approach and the subchannel wise Q-WSRM scheme using the WMMSE design [6]. The metric used for the comparison is the total number of backlogged bits left over after each transmission, which is denoted as  $\chi = \sum_{k=1}^K [Q_k - t_k]^+$ . Even though  $\mathcal{U}(1)$  and  $\mathcal{U}(3)$  has equal number of backlogged packets of  $Q_1 = Q_3 = 4$  bits, user  $\mathcal{U}(3)$  is scheduled in the first sub-channel due to the better channel condition. In contrast, the JSFRA approach assigns the first user on the first sub-channel, which reduces the total number of backlogged packets. The rate allocated for  $\mathcal{U}(2)$  on the second sub-channel is higher in JSFRA scheme compared to the others. It is due to the efficient allocation of the total power shared across the sub-channels.

For a MIMO scenario, we consider a system with N=3 sub-channels and  $N_B=3$  BSs, each equipped with  $N_T=4$  transmit antennas operating at 10dB SNR, serving  $|\mathcal{U}_b|=3$  users each. The PL between the BSs and the users are uniformly generated from [0,-3] dB and the associations are made by selecting the BS with the lowest PL component. Fig. 1(a) shows the performance of the centralized schemes for a single receive antenna system. The total number of queued packets for Fig. 1(a) is given by  $Q_k=[14,15,14,8,12,9,12,11,11]$  bits and for Fig. 1(b) is  $Q_k=[9,12,8,12,5,4,10,8,5]$  bits respectively.

The comparison is made in terms of the total number of residual bits remaining in the system after each SCA update in Fig. 1. The Q-WSRM algorithm is not optimal due to the problem of over-allocation when the number of queued packets are few in number. In contrast, the Q-WSRME algorithm provides more favorable allocation with the explicit rate constraint to avoid the over-allocation. For both scenarios in Fig. 1, the Q-WSRME performs marginally inferior to the JSFRA algorithms due to the weights used in the algorithm. It is due to the fact that the Q-WSRME algorithm favors the users with the large number of backlogged packets as compared to the users with better channel conditions. Fig. 1(b) compares the algorithms for  $N_R = 2$  receive antenna case. In all figures, the receivers are updated along with the SCA update instants i.e,  $J_{\text{max}} = 1$  in Algorithm 1. The degradation by performing combined update is marginal, since the receiver minimizes the objective for a fixed transmit precoders.

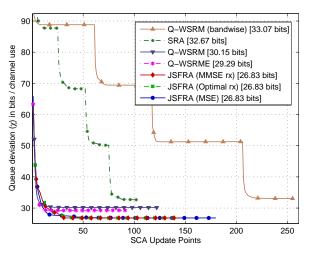
The behavior of the JSFRA algorithm for different exponents q is outlined in the Table II for the users located at the cell-edge of the system employing  $N_T=4$  transmit antennas. It is evident that the JSFRA algorithm minimizes the total number of queued bits for the  $\ell_1$  norm compared to the  $\ell_2$  norm, which is shown in the column displaying the total number of left over packets  $\chi$  in bits. The  $\ell_\infty$  norm provides fair allocation of the resources by making the left over packets to be equal for all users to  $\chi_k=3.58$  bits.

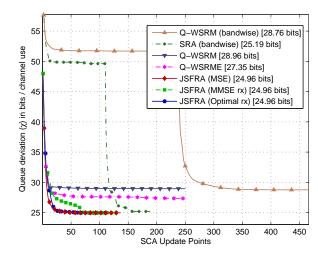
#### B. Distributed Solutions

The distributed algorithms are compared using the total number of backlogged packets after each SCA update. Fig. 2 compares the performances of the algorithms with the PL varies uniformly between [0,-6] dB. Each BS serves  $|\mathcal{U}_b|=4$  users in a coordinated manner to reduce the number of backlogged packets, whith the queued packets for each user is given by  $Q_k=[5,7,9,11,8,12,5,4]$  bits. As discussed in

TABLE I
SUB-CHANNEL-WISE LISTING OF CHANNEL GAINS AND RATE ALLOCATIONS BY DIFFERENT ALGORITHMS FOR A SCHEDULING INSTANT

Users	Queued Packets	Channel Gains			Q-WSRME approach (modified <i>backpressure</i> )			JSFRA Scheme			Q-WSRM band Alloc Scheme		
		SC-1	SC-2	SC-3	SC-1	SC-2	SC-3	SC-1	SC-2	SC-3	SC-1	SC-2	SC-3
1	4	1.71	0.53	0.56	0	0	0	4.0	0	0	0	0	0
2	8	0.39	1.41	1.03	0	4.88	3.11	0	5.49	0	0	4.39	3.53
3	4	2.34	1.26	2.32	4.0	0	0	0	0	4.0	5.81	0	0
Re	Remaining backlogged packets (χ)				3.92 bits			2.51 bits			5.89 bits		





(a). System Model  $\{N,N_B,K,N_T,N_R\}=\{4,3,9,4,1\}$ 

(b). System Model  $\{N, N_B, K, N_T, N_R\} = \{2, 3, 9, 4, 2\}$ 

Fig. 1. Total number of backlogged packets  $\chi$  present in the system after each SCA updates using  $\ell_1(q=1)$  norm for JSFRA schemes

TABLE II Number of backlogged bits associated with each user for a system  $\{N,N_B,K,N_R\}=\{5,2,8,1\}.$ 

_ a	user indices									
q	1	2	3	4	5	6	7	8	χ	
1	15.0	3.95	5.26	8.95	7.0	11.9	12.0	9.7	25.15	
2	11.2	3.9	10.76	10.65	10.27	9.68	8.77	5.9	27.77	
$\infty$	11.4	4.4	10.4	10.4	10.4	8.4	8.4	6.4	28.68	
$Q_k$	15.0	8.0	14.0	14.0	14.0	12.0	12.0	10.0		

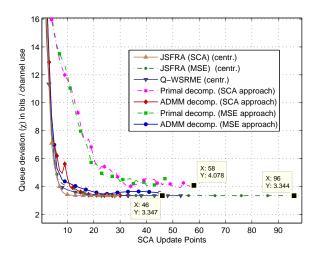


Fig. 2. Convergence of the centralized and the distributed algorithms for  $\{N, N_B, K, N_T, N_R\} = \{3, 2, 8, 4, 1\}$  using  $\ell_1$  norm for JSFRA schemes

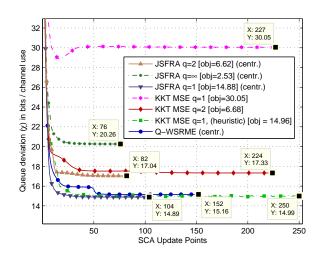


Fig. 3. Impact of varying q in the total number of backlogged packets after each SCA update for a system  $\{N, N_B, K, N_T, N_R\} = \{5, 2, 8, 4, 1\}$ 

Section IV, the performance and the convergence speed of the distributed algorithms are susceptible to the step size used in the subgradient update. Due to the fixed interference levels in the primal approach, it may lead to infeasible solutions if the initial or any intermediate update is not feasible.

Fig. 2 compares the primal and the ADMM solutions for the JSFRA scheme using the SCA and by MSE relaxation using the number of backlogged packets. In between the SCA updates, the primal or the ADMM scheme is performed for  $J_{\rm max}=20$  iterations to exchange the respective coupling variables. The number of backlogged packets only at the SCA points are marked in the figure. The performance of the distributed approaches are similar to the centralized schemes if the primal and dual updates are performed until convergence.

Fig. 3 compares the performances of the centralized and the KKT algorithm in Section IV-C for different exponents with  $Q_k = [9, 16, 14, 16, 9, 13, 11, 12]$  bits and the PL varies uniformly between [0, -3] dB. The  $\ell_1$  norm JSFRA scheme performs better over other schemes due to the greedy objective. The KKT approach for  $\ell_1$  norm is not defined due to the nondifferentiability of the objective as discussed in the Section IV-C. If used for  $\ell_1$  norm, the over-allocation will not affect the dual variables  $\sigma_{l,k,n}$  and  $\alpha_{l,k,n}$  since the queue deviation is raised to the power zero in (42e). A heuristic method is proposed in Fig. 3 by assigning zero for  $\sigma_{l,k,n}$  when  $Q_k - t_k < 0$  to addresses the over-allocation. The heuristic approach oscillates near the converging point with the deviation determined by the factor  $\rho$  used in (42f). The objective values are mentioned in the legend for all the schemes and the  $\ell_1$  norm is used for comparison.

#### C. Queuing Analysis over Multiple Transmission Slots

In this section, we numerically study the performance of the centralized algorithms with different  $\ell_q$  values over multiple transmission slots. The system model examined for the illustrations is provided in Fig. 4. For all users in the system, the average arrivals  $A_k$ 's are fixed and varied equally for the model considered in Fig. 4(a), and for Fig. 4(b), the average arrival is fixed to be  $A_k=6$  bits. Note that the instantaneous arrivals  $\lambda_k(i)$  are all different and it follows the Poisson process. The PL is modeled as a uniform random variable [0,-6] dB.

Fig. 4(a) plots the average of the total number of backlogged packets left out in the system after each transmission instant, i.e,  $E_i\left[\sum_k\left[Q_k(i)-t_k(i)\right]^+\right]$ . Unlike the Q-WSRM scheme, the average backlogged packets of the  $\ell_2$  JSFRA scheme is comparable to the Q-WSRME approach for all average arrival rates due to the explicit rate constraints (14). However, when  $A_k \geq 7$  bits in Fig. 4(a), both Q-WSRM and Q-WSRME schemes perform the same since the problem of over-allocation is negligible. The performance of the  $\ell_1$  JSFRA scheme outperforms all other schemes in terms of the average number of residual packets due to the greedy allocation at each instant.

Fig. 4(a) also includes the uncoordinated  $\ell_1$  JSFRA scheme and the time division multiplexing (TDM) mode without considering the inter-cell interference terms in the SINR expressions while designing the precoders. The TDM scheme with the fixed total power constraint performs inferior to the uncoordinated transmission due to the diverse user PL variations considered. Fig. 4(b) compares the total number of backlogged packets left in the system after each transmission slot by different centralized algorithms. The total number of residual packets for the Q-WSRM scheme is noticeably large in comparison with the other schemes in Fig. 4(b) due to the inability in controlling the over-allocations at each instant. The instantaneous fairness constraint imposed by the  $\ell_{\infty}$  JSFRA

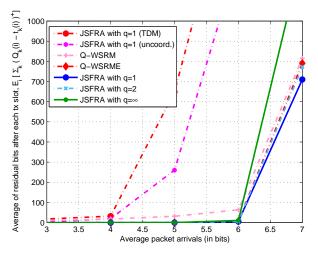
scheme is effective in reducing the number of backlogged packets for the average arrival rate considered in Fig. 4(b). However, in Fig. 4(a), the  $\ell_{\infty}$  JSFRA performs inferior to the Q-WSRM scheme, since the fairness is not effective when the system is unstable, *i.e*, when  $A_k \geq 7$  in the current model.

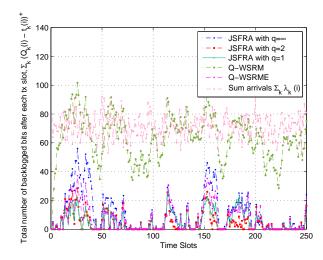
#### VI. CONCLUSIONS

In this paper, we addressed the problem of allocating downlink space-frequency resources to the users in a multi-cell MIMO IBC system using OFDM. The resource allocation is considered as a joint space-frequency precoder design problem since the allocation of a resource to a user is obtained by a non-zero precoding vector. We proposed the JSFRA scheme by relaxing the nonconvex DC constraint by a sequence of convex subsets using the SCA for designing the precoders to minimize the total number of user queued packets. Additionally, an alternative MSE reformulation approach is also proposed by using the SCA to address the nonconvex DC constraints for a fixed MMSE receivers. We also proposed various methods to decentralize the precoder designs for the JSFRA problem using primal and ADMM methods. Finally, we proposed a practical iterative algorithm to obtain the precoders in a decentralized manner by solving the KKT conditions of the MSE reformulated JSFRA method. The proposed iterative algorithm requires few iterations and limited signaling exchange between the coordinating BSs to obtain the efficient precoders for a given number of iterations. Numerical results are used to compare the performances of the proposed algorithms.

#### REFERENCES

- E. Matskani, N. Sidiropoulos, Z.-Q. Luo, and L. Tassiulas, "Convex Approximation Techniques for Joint Multiuser Downlink Beamforming and Admission Control," *IEEE Trans. Wireless Commun.*, vol. 7, no. 7, pp. 2682–2693, July 2008.
- [2] C. Ng and H. Huang, "Linear Precoding in Cooperative MIMO Cellular Networks with Limited Coordination Clusters," *IEEE J. Sel. Areas Commun.*, vol. 28, no. 9, pp. 1446–1454, December 2010.
- [3] L.-N. Tran, M. Hanif, A. Tölli, and M. Juntti, "Fast Converging Algorithm for Weighted Sum Rate Maximization in Multicell MISO Downlink," *IEEE Signal Process. Lett.*, vol. 19, no. 12, pp. 872–875, 2012
- [4] 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, Nov 2007.
- [5] S. S. Christensen, R. Agarwal, E. Carvalho, and J. Cioffi, "Weighted Sum-Rate Maximization using Weighted MMSE for MIMO-BC Beamforming Design," *IEEE Trans. Wireless Commun.*, vol. 7, no. 12, pp. 4792–4799, 2008.
- [6] Q. Shi, M. Razaviyayn, Z.-Q. Luo, and C. He, "An Iteratively Weighted MMSE Approach to Distributed Sum-Utility Maximization for a MIMO Interfering Broadcast Channel," *IEEE Trans. Signal Process.*, vol. 59, no. 9, pp. 4331–4340, sept. 2011.
- [7] M. Hong, Q. Li, Y.-F. Liu, and Z.-Q. Luo, "Decomposition by Successive Convex Approximation: A Unifying Approach for Linear Transceiver Design in Interfering Heterogeneous Networks," 2012. [Online]. Available: http://arxiv.org/abs/1210.1507
- [8] J. Kaleva, A. Tölli, and M. Juntti, "Primal Decomposition based Decentralized Weighted Sum Rate Maximization with QoS Constraints for Interfering Broadcast Channel," in *IEEE 14th Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, 2013, pp. 16–20.
- [9] ——, "Decentralized Beamforming for Weighted Sum Rate Maximization with Rate Constraints," in 24th International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC Workshops). IEEE, 2013, pp. 220–224.





- (a). Average backlogged packets in the system after 250 transmission instants
- (b). Total backlogged packets at each transmission slot for  $A_k = 6$  bits

Fig. 4. Time analysis of the Queue dynamics for a system  $\{N, N_B, K, N_T, N_R\} = \{4, 2, 12, 4, 1\}$ 

- [10] D. P. Palomar and M. Chiang, "A Tutorial on Decomposition Methods for Network Utility Maximization," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 8, pp. 1439–1451, 2006.
- [11] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers," *Foundations and Trends* (in Machine Learning, vol. 3, no. 1, pp. 1–122, 2011.
- [12] H. Pennanen, A. Tölli, and M. Latva-Aho, "Decentralized Coordinated Downlink Beamforming via Primal Decomposition," *IEEE Signal Pro*cess. Lett., vol. 18, no. 11, pp. 647–650, 2011.
- [13] A. Tölli, H. Pennanen, and P. Komulainen, "Decentralized Minimum Power Multi-Cell Beamforming with Limited Backhaul Signaling," *IEEE Trans. Wireless Commun.*, vol. 10, no. 2, pp. 570–580, 2011.
- [14] L. Tassiulas and A. Ephremides, "Stability Properties of Constrained Queueing Systems and Scheduling Policies for Maximum Throughput in Multihop Radio Networks," *IEEE Trans. Autom. Control*, vol. 37, no. 12, pp. 1936–1948, Dec 1992.
- [15] M. Neely, Stochastic Network Optimization with Application to Communication and Queueing Systems, ser. Synthesis Lectures on Communication Networks. Morgan & Claypool Publishers, 2010, vol. 3, no. 1.
- [16] L. Georgiadis, M. J. Neely, and L. Tassiulas, Resource Allocation and Cross-Layer Control in Wireless Networks. Now Publishers Inc, 2006.
- [17] R. A. Berry and E. M. Yeh, "Cross-Layer Wireless Resource Allocation," IEEE Signal Process. Mag., vol. 21, no. 5, pp. 59–68, 2004.
- [18] M. Chiang, S. Low, A. Calderbank, and J. Doyle, "Layering as Optimization Decomposition: A Mathematical Theory of Network Architectures," *Proceedings of the IEEE*, vol. 95, no. 1, pp. 255–312, Jan 2007.
- [19] K. Seong, R. Narasimhan, and J. Cioffi, "Queue Proportional Scheduling via Geometric Programming in Fading Broadcast Channels," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 8, pp. 1593–1602, 2006.
- [20] P. C. Weeraddana, M. Codreanu, M. Latva-aho, and A. Ephremides, "Resource Allocation for Cross-Layer Utility Maximization in Wireless Networks," *IEEE Trans. Veh. Technol.*, vol. 60, no. 6, pp. 2790–2809, 2011.
- [21] F. Zhang and V. Lau, "Cross-Layer MIMO Transceiver Optimization for Multimedia Streaming in Interference Networks," *IEEE Trans. Signal Process.*, vol. 62, no. 5, pp. 1235–1244, March 2014.
- [22] B. R. Marks and G. P. Wright, "A General Inner Approximation Algorithm for Nonconvex Mathematical Programs," *Operations Research*, vol. 26, no. 4, pp. 681–683, 1978.
- [23] Z.-Q. Luo and S. Zhang, "Dynamic Spectrum Management: Complexity and Duality," *IEEE J. Sel. Topics Signal Process.*, vol. 2, no. 1, pp. 57– 73, Feb 2008.
- [24] S. P. Boyd and L. Vandenberghe, Convex optimization. Cambridge University Press, 2004. [Online]. Available: http://stanford.edu/~boyd/ cvxbook/bv\_cvxbook.pdf
- [25] T. Lipp and S. Boyd, "Variations and Extensions of the Convex-Concave Procedure," 2014. [Online]. Available: http://web.stanford.edu/~boyd/ papers/pdf/cvx\_ccv.pdf

- [26] G. R. Lanckriet and B. K. Sriperumbudur, "On the Convergence of the Concave-Convex Procedure," in Advances in Neural Information Processing Systems, 2009, pp. 1759–1767.
- [27] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.0 beta," http://cvxr.com/cvx, Sep. 2013.
- [28] P. Komulainen, A. Tölli, and M. Juntti, "Effective CSI Signaling and Decentralized Beam Coordination in TDD Multi-Cell MIMO Systems," *IEEE Trans. Signal Process.*, vol. 61, no. 9, pp. 2204–2218, 2013.
- [29] C. Shen, T.-H. Chang, K.-Y. Wang, Z. Qiu, and C.-Y. Chi, "Distributed Robust Multicell Coordinated Beamforming With Imperfect CSI: An ADMM Approach," *IEEE Trans. Signal Process.*, vol. 60, no. 6, pp. 2988–3003, June 2012.
- [30] R. Meyer, "Sufficient Conditions for the Convergence of Monotonic Mathematical Programming Algorithms," *Journal of Computer and System Sciences*, vol. 12, no. 1, pp. 108–121, 1976.
- [31] W. Rudin, Principles of Mathematical Analysis. McGraw-Hill New York, 1964, vol. 3.
- [32] G. Scutari, F. Facchinei, L. Lampariello, and P. Song, "Distributed Methods for Constrained Nonconvex Multi-Agent Optimization – Part I: Theory." [Online]. Available: http://arxiv.org/abs/1410.4754v1
- [33] J. C. Bezdek and R. J. Hathaway, "Some Notes on Alternating Optimization," in *Advances in Soft Computing AFSS*. Springer, 2002, pp. 288–300.
- [34] T. D. Quoc and M. Diehl, "Sequential Convex Programming Methods for Solving Nonlinear Optimization Problems with DC Constraints," 2011. [Online]. Available: http://arxiv.org/abs/1107.5841v1
- [35] W. Zangwill, Nonlinear Programming: A Unified Approach, ser. Prentice-Hall International Series in Management. Prentice-Hall, 1969.
- [36] D. P. Bertsekas, Nonlinear Programming, 2nd ed. Athena Scientific, sep 1999.
- [37] D. P. Bertsekas and J. N. Tsitsiklis, Parallel and Distributed Computation: Numerical Methods. Prentice Hall Englewood Cliffs, NJ, 1989, vol. 23