# Adaptive utility-based scheduling algorithm for multiuser MIMO uplink

Tine Celcer[1*], Gorazd Kandus[2] and Tomaž Javornik[2]

## Abstract

Resource allocation issues are discussed in the context of a virtual multiuser MIMO uplink assuming users equipped with a single antenna. A scheduling algorithm, which efficiently mitigates the co-channel interference (CCI) arising from the spatial correlation of users sharing common resources, is proposed. Users are selected using an incremental approach with a reduced complexity that is due to the elimination of over-correlated users at each iteration. The user selection criterion is based on an adaptive, utility-based scheduling metric designed for the purpose. Its main advantage lies in the periodic adaptation of priority weights according to the application characteristics described with its utility curves and according to momentary quality of service (QoS) parameters. The results show a better performance in aggregate system utility than the existing utility based scheduling metrics such as proportionally fair scheduling (PFS), largest weighted delay first (LWDF), modified LWDF (M-LWDF), and exponential algorithm.

**Keywords:** Multiuser systems, Adaptive resource allocation, Utility, MIMO, ACM

## Introduction

Over the last two decades, achievements in the field of transmission techniques have enabled the transmission of data with high throughput in wireless systems [1,2]. The area of wireless communication networks and technologies has evolved and is still evolving at a high pace [3]. One of the consequences is a wide range of applications supported by user terminals and services provided by network operators. Heterogeneous classes of service requiring high reliability of transmission and/or high throughput, along with low transmission delays, make the provision of quality of service (QoS) in wireless systems a challenging task, due to the scarcity of wireless resources. As the bandwidth and transmission power are limited resources, it is important to exploit the given spectrum effectively in order to maximize the number of users achieving the desired QoS level.

Among other advances, a significant increase in throughput and/or transmission reliability may be achieved by using multiple antennas at the receiver and transmitter ends, thus enabling efficient exploitation of physical wireless channel properties in the spatial domain [2]. The so-called multiple input multiple output (MIMO) systems take advantage of the multipath signal spreading, considered as a detrimental characteristic of the wireless channel in single antenna systems. The increase in throughput, of an order equal to a minimum number of transmit and receive antennas, can be achieved by multiplexing independent data streams across different transmit antennas with the application of a V-BLAST transmission scheme [4]. However, mobile terminals are usually equipped only with a single antenna, which prevents the use of this technique on a point-to-point link, since pursuant to the theory of spatial multiplexing, the number of receive antennas has to be equal to or higher than the number of simultaneously transmitted independent data streams [5]. Nevertheless, even in such cases, spatial multiplexing of user streams may be applied in multiuser systems by way of using a spatial domain multiple access (SDMA) scheme. The base station (BS) equipped with multiple antennas and users equipped with a single antenna and sharing common radio resources are thus forming a virtual MIMO system. Due to this virtuality, a fundamental difference between uplink and downlink user grouping process exists.

* Correspondence: tine.celcer@cobik.si
[1]The Centre of Excellence for Biosensors, Instrumentation and Process Control - COBIK, Velika pot 22, SI-5250 Solkan, Slovenia
Full list of author information is available at the end of the article

In general, there are no direct communication links between users, hence the cooperation between users is not possible in the downlink, and the approaches known from single link MIMO systems cannot be applied directly. However, an appropriate precoding technique, responsible for inter-user interference mitigation, may be applied at the transmitter to make spatial user grouping possible. Examples of such user grouping methods are theoretically optimal dirty paper coding (DPC) [6] and various less complex but suboptimal beam-forming techniques [7-9].

Complex precoding techniques are not required in the uplink due to sufficient processing capabilities at the BS. Nevertheless, the absence of user grouping precoding techniques reflects in co-channel interference (CCI) due to the correlation of spatial signatures of users sharing common radio resources. In order to mitigate the CCI effectively and provide high system level efficiency, a set of spatially multiplexed users has to be selected carefully, making user scheduling one of the most crucial areas of resource management. Resource allocation algorithms with scheduling metrics, based on utility optimization, have proved to be strong candidates for solving the resource allocation problem, since their major advantage lies in strong coupling between user satisfaction and system level efficiency [10].

Based on the type of the parameters considered for utility definition, the existing utility-based scheduling metrics can be divided into three groups, namely, throughput maximization oriented channel-aware algorithms, delay optimization queue-aware algorithms and channel- and queue-aware scheduling algorithms that combine the parameters from different layers of the protocol stack.

Throughput maximization oriented algorithms, i.e. maximal rate and proportional fair scheduling (PFS) algorithms [11], with channel-dependent scheduling metrics yield high aggregate throughput by exploiting multiuser diversity [12]. However, they only perform well in networks with homogeneous, delay-tolerant traffic and with a sufficient level of user mobility. In the case delay-sensitive, real-time (RT) traffic is present, they cannot satisfy diverse QoS requirements, since they prioritize users with good channel conditions without considering packet waiting time and traffic priority. Therefore, the system level efficiency in networks with heterogeneous traffic should not only be characterized by aggregate system throughput but also, and most importantly, by QoS level and satisfaction of each user.

The Largest Weighted Delay First (LWDF) scheduling algorithm [13], on the other hand, provides QoS differentiation for RT traffic by considering the current delay of packets in the queue, weighted with a traffic priority factor. However, the LWDF rule disregards any kind of

channel state information (CSI), thus preventing the exploitation of time-varying link conditions.

In order to optimize the system level efficiency, it is important that a scheduling metric combines QoS related parameters (packet waiting time and priority weights, depending on the class of service) with channel-dependent information. Pursuing this objective, the so-called throughput-optimal scheduling algorithms, such as the Modified-Largest Weighted Delay First (M-LWDF) rule [14] and the exponential (EXP) rule [15], improve the quality of resource allocation significantly. Throughput optimal policy is defined as a policy that can keep the queues stable for all users in the system, providing this is at all made feasible with any of the scheduling policies.

Nevertheless, throughput optimality does not explicitly guarantee the provision of QoS in the form of delay or throughput bounds, and different throughput-optimal algorithms show different performance or fairness properties. Hence, there is still potential for further improvement in scheduling algorithm design. In the light of the aforementioned, certain drawbacks of M-LWDF and EXP algorithms can be identified. First, their metrics do not consider the different shapes of the utility curves as a function of throughput or packet delay as per different classes of service, and secondly, the priority weights are constants calculated on the basis of the statistical definition of QoS requirements, expressed in terms of the probability of maximal packet delay violation. Consideration of the utility curves and their characteristics, in combination with periodic priority weight adaptation, can further increase the system level efficiency.

In this article, we propose a novel scheduling algorithm with an adaptive, utility-based scheduling metric for the multiuser MIMO uplink, together with the support for SDMA. The study is limited to the case where users are equipped with single antenna terminals. The CCI is mitigated efficiently using a maximal correlation threshold for users sharing common resources, while the scheduling metric is derived from the M-LWDF scheduling rule, with the main difference being that the static priority weights are substituted by adaptive weights, thus increasing the flexibility of the scheduling metric according to instantaneous system requirements. Adaptation of the priority weights is performed based on the ratio between the momentary and the target values of QoS parameters for different traffic types. The algorithm also enables the selection of optimal transmission modes for selected users by using a linear zero-forcing (ZF) detection algorithm at the receiver, since the SNR, achieved after detection, can be analytically calculated in advance.

The remainder of the article is organized as follows. In 'Utility curves for different types of traffic' section, the

performance characteristics of different traffic types as a function of packet delay and allocated bandwidth are presented. Next, we describe the design of the proposed scheduling algorithm, with an emphasis on the adaptation of priority weights. In 'Wireless system model and algorithm parameters' section, the system model and algorithm parameters are presented, while the algorithm performance evaluation is given in 'Performance analysis' section. Conclusions are drawn in 'Conclusion' section.

### Utility curves for different types of traffic

Normalized packet utility, in terms of the allocated bandwidth (or, equally, transmission rate), is depicted in Figure 1[16]. The utility curve for delay-tolerant, best-effort (BE) data traffic is characterized by a monotonically increasing function, with decreasing marginal improvement as the packet transmission rate increases (Figure 1a). The elastic nature of such applications is characterized by a strong adaptivity to delay and bandwidth. Hard RT applications, such as VoIP, have a utility function with the shape of a step function (Figure 1b). These applications require the packets to be transmitted inside a given delay bound. If the packet arrives too late (i.e. the transmission rate is on average lower than the data arrival rate), it proves useless, and the user satisfaction level, i.e. packet utility, equals zero. When the threshold is achieved, user satisfaction level increases instantly, and no further increase is achieved with an additional bandwidth allocation (higher transmission rate). Due to the possibility of adjusting their data generation rate through scalable coding some RT applications, such as video streaming, have a certain level of adaptivity to delay and allocated transmission rate. Their
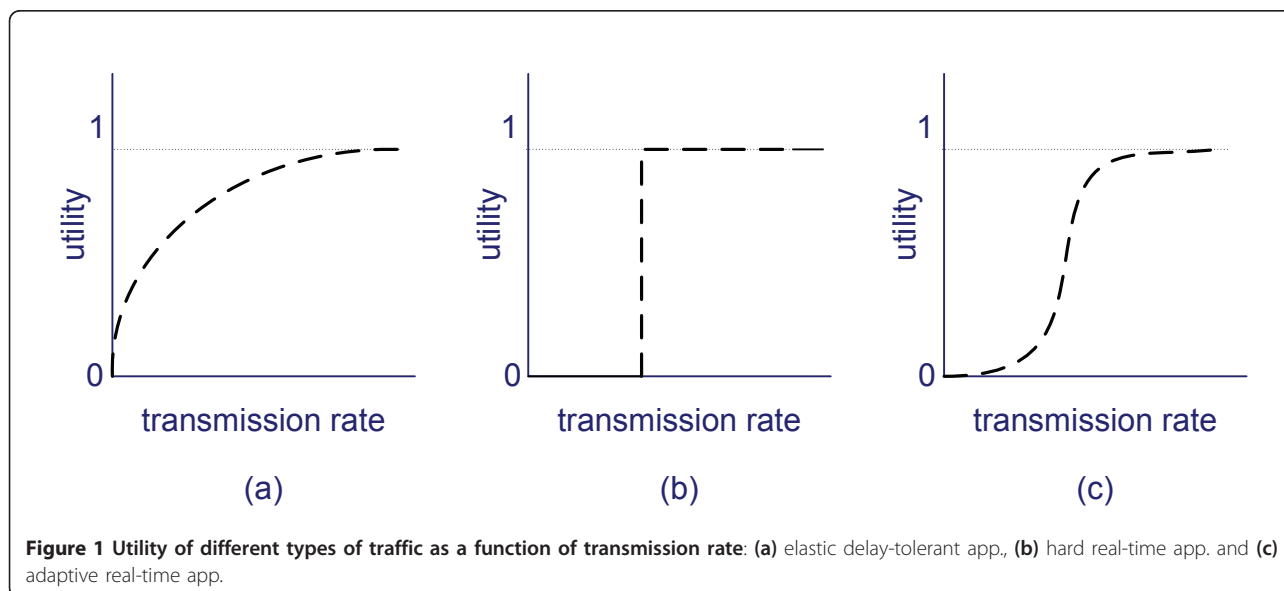
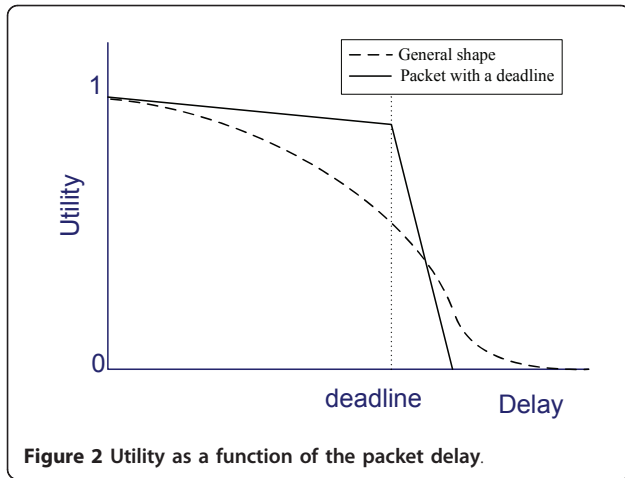utility curve is smoother than that of the hard RT applications (Figure 1c).

The aforementioned characteristics of the different traffic types show why it proves important to take such features into consideration in the design of scheduling metric. The impact of an equal decrease in the allocated transmission rate on packet utility, i.e. user satisfaction level, is not the same for the RT user as it is for the BE user. Disregarding this fact will significantly influence the aggregate system efficiency.

The utility of transmitted packets for delay-sensitive applications can also be presented as a function of packet end-to-end delay, consisting of packet queuing delay and transmission delay. Corresponding normalized utility curves are presented in Figure 2[17]. In this case, the utility is a monotonically decreasing function, presenting an incremental marginal decrease as the delay increases. In general, the utility has a smooth form (dashed line); however, if the packet has a deadline, the utility (solid line) is relatively flat (the application disregards if the packet arrives earlier), and drops sharply when the deadline (vertical dotted line) is passed.

### Proposed adaptive scheduling algorithm with SDMA support

In this section, the design of a cross-layer scheduling algorithm for networks with heterogeneous traffic types is presented. The algorithm can be divided into three mutually dependent steps (Figure 3), namely, CCI mitigation and user grouping (blue coloured blocks with a solid line), user selection, based on the proposed adaptive scheduling metric using an incremental approach (green coloured blocks with a dashed line) and optimal



**Figure 1 Utility of different types of traffic as a function of transmission rate**: **(a)** elastic delay-tolerant app., **(b)** hard real-time app. and **(c)** adaptive real-time app.

**Figure 2 Utility as a function of the packet delay**.

transmission scheme selection (yellow coloured block with a dashed-dotted line).

The algorithm is designed for a single cell, multi-user distributed MIMO system, where the base station (BS) is equipped with $M$ antennas serving $K$ active users, each equipped with a single antenna. In general, the proposed algorithm can be applied for both downlink and uplink; however, in this article, the study is limited to uplink communication only, where additional pre-coding is not required, as explained in the 'Introduction' section.

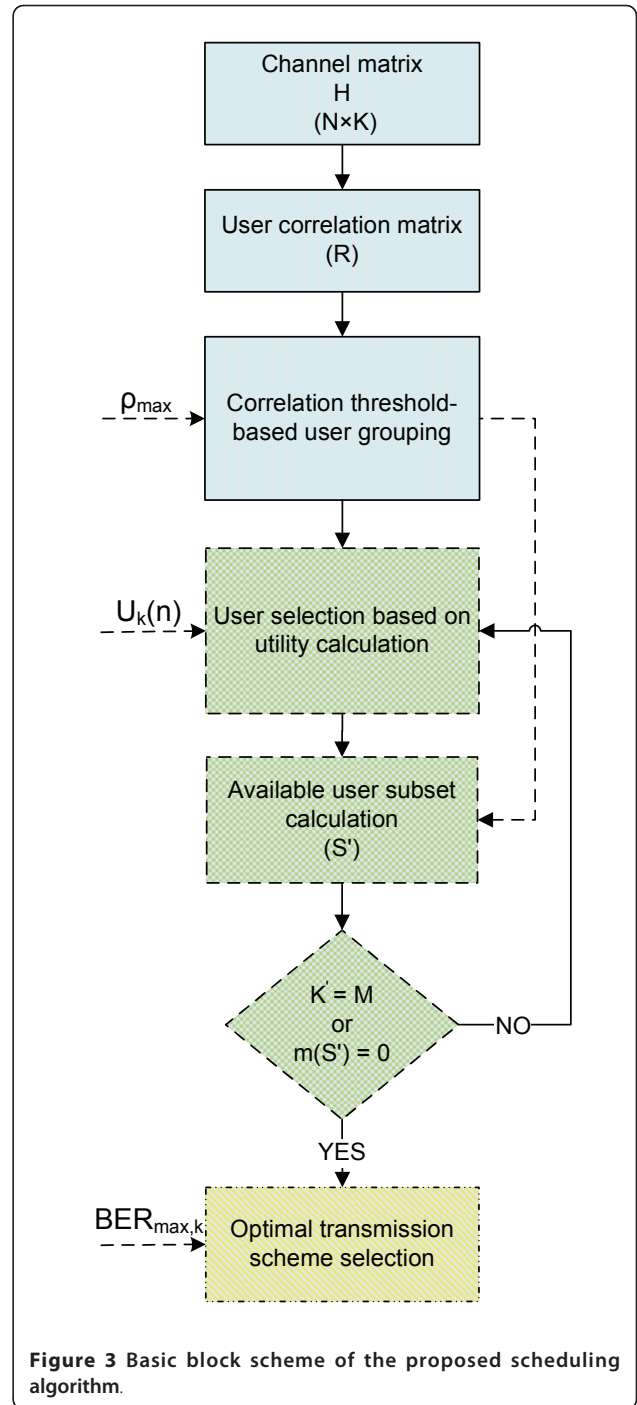*User grouping and CCI mitigation*

To separate spatially multiplexed data streams, the use of a linear ZF receiver is assumed, mainly due to its simplicity and low computational complexity. However, linear ZF receivers suffer from noise enhancement, especially, if the user spatial signatures are highly correlated; it is crucial, therefore, to limit the CCI. For that reason, the algorithm first calculates the correlation matrix $\mathbf{R}$, using the channel matrix $\mathbf{H}$, which can be used to describe frequency flat fading MIMO systems [2,5] and is composed of the users' $\mathbf{M} \times 1$ channel vectors $\mathbf{h}_k$. First, each channel vector $\mathbf{h}_k$ is normalized, so that $\|\mathbf{h_k}\|_F^2 = 1$:

$$\mathbf{h_{k\_norm}} = \frac{\mathbf{h_k}}{\sqrt{\left|\mathbf{h_k}^*\mathbf{h_k}\right|}}.$$

(1)

Matrix $\mathbf{R}$ is then calculated, using the equation:

$$\mathbf{R} = \left|\mathbf{H_{norm}^*} \cdot \mathbf{H_{norm}}\right| = \begin{bmatrix} 1 & \rho_{12} & \cdots & & \rho_{1K} \\ \rho_{12} & 1 & \ddots & & \vdots \\ \vdots & \ddots & 1 & \rho_{(K-1)K} \\ \rho_{K1} & \cdots & \rho_{K(K-1)} & 1 \end{bmatrix},$$

(2)

where $\mathbf{H}_{norm}$ is composed of normalized channel vectors $\mathbf{h}_{k\_norm}$. The elements $\rho_{ij}$ $(i,j = 1,...,K)$ represent the correlation between the $i$th and $j$th users.



**Figure 3 Basic block scheme of the proposed scheduling algorithm**.

CCI is mitigated with the introduction of the maximal allowed correlation between any pair of users sharing the same resources ($\rho_{\max}$). By adopting this approach, the CCI can be mitigated to an arbitrary level. Next, a group of users $\mathbf{S}_k$ meeting the following condition is defined for each user:

$$\mathbf{S}_k = \left\{j; j \neq k, \rho_{jk} < \rho_{\max}\right\} \; ; \; k = 1, \ldots, K.$$

(3)

The group $\mathbf{S}_k$ thus contains all those users allowed to share common resources with user $k$. Note that each user can be placed in a number of groups. This approach is based on the idea presented in [18], where the authors propose to form several groups of users, based on the maximal allowed correlation. The users in the same group cannot share channel resources simultaneously, while the correlation between any pair of users from different groups is lower than the threshold value. The proposed grouping is complicated and leads to inadequate situations. The user grouping, proposed in this article, eliminates this deficiency.

When users are grouped, the incremental approach is used to select a set of spatially multiplexed users. In each iteration, the radio resources are allocated to the user with the highest metric among all active users. The novel, adaptive utility-based scheduling metric, is explained in detail in the next subsection. We start with a full set of active users and, after each iteration, update the set of available users $\mathbf{S}'$ by eliminating over-correlated users. If the $k$th user is chosen, then $\mathbf{S}'$ for the $i$th iteration is updated as follows:

$$\mathbf{S}'(i) = \mathbf{S}'(i-1) \cap \mathbf{S}_k. \qquad (4)$$

We repeat the iterations as long as the number of selected users is smaller or equal to the number of receive antennas at the BS, or as long as $m(\mathbf{S}') > 0$.

The advantage of this approach is twofold. First, the interference is limited in a simple and effective way, thus keeping the scheduling metric simple, since no parameter based on any relation between users is required, and the utility does not have to be recalculated after each iteration. Secondly, the complexity of user selection is decreased, since the search space is reduced after each iteration. The reduction of the search space in the case of $M = 4$ and $\rho_{\max} = 0.5$, averaged over 20,000 independent channel realizations, is depicted in Figure 4, where (a) indicates the number of available users in different iterations, and (b) the ratio between the number of available users and the full set of users. In the case of the basic incremental approach, i.e. $\rho_{\max} = 1$, the number of available users in the $i$th iteration is $K -(i - 1)$. Simulations have shown that the cardinality of $\mathbf{S}'$ is decreased from around 50% after the first iteration, to less than 30% after the second one and, down to only around 10% of the full set after the third iteration. Naturally, the advantage of such an approach is evident in the case of a large number of users, where the level of multiuser diversity is high and 'good' users may be found even if the search space is significantly reduced. Moreover, the reduction of the search space depends on the selected value of the parameter $\rho_{\max}$. The

optimization of this parameter will be presented in 'Wireless system model and algorithm parameters' section.

### Utility-based scheduling metric

In each iteration, the decision on the user selection is made by using a channel-and queue-aware scheduling metric, derived from the M-LWDF approach [14]. The drawback of the M-LWDF scheduling algorithm, when deployed in a heterogeneous service scenario, is its characteristic to maintain the stability of the queues, and this does not necessarily guarantee low delays. BE traffic might occupy the bandwidth and consequently insufficient amount of resources is assigned to RT traffic, preventing the provision of required QoS levels. The adaptation of M-LWDF approach to a mixed service scenario has also been investigated in [19] by manipulating $T_i$ and $\delta_i$ parameters. The main advantage of the scheduling metric, proposed in this article, is the adaptivity of its priority weights, taking into consideration the specific shapes of the utility curves, as presented in 'Utility curves for different types of traffic' section. The real-time tuning of the priority weights is based on the ratio between the actual and target values of the QoS parameters, namely, transmission rate and maximal delay.
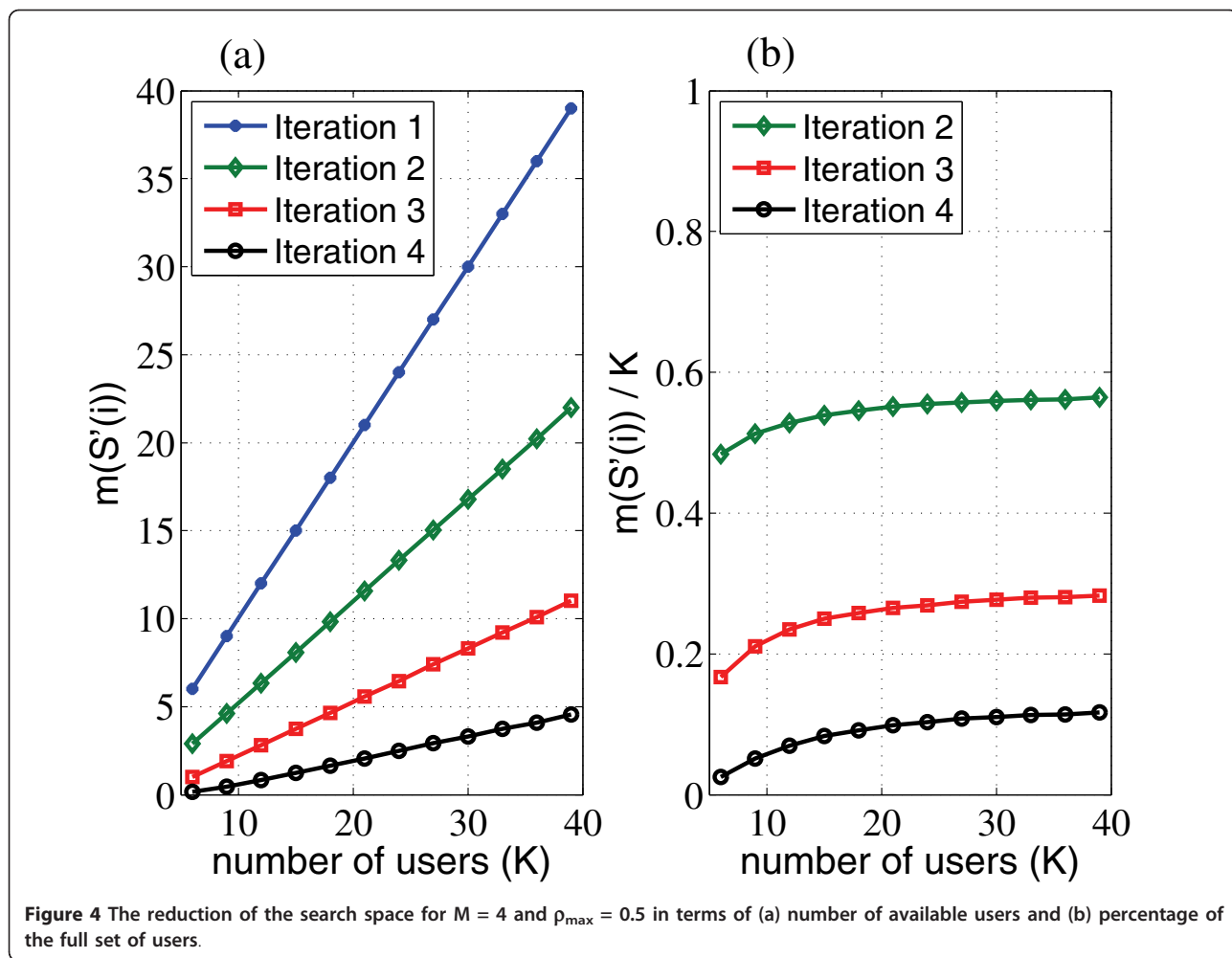
In the proposed algorithm, the utility for the $k$th user in the $n$th time frame is calculated using the following scheduling metric:

$$U_k(n) = d_{\mathrm{HOL},k}(n)^{a_k(n)} \cdot \frac{r_k(n)}{\bar{r}_k} \cdot b_k(n), \qquad (5)$$

where $d_{\mathrm{HOL},k}(n)$ is the waiting time of the head-of-line (HOL) packet, $r_k(n)$ the theoretically achievable transmission rate in an interference free environment, and $\bar{r}_k$ the average transmission rate. The utility function introduces two adaptive weights, i.e. a delay-dependent weight $a_k(n)$ and a throughput-dependent weight $b_k(n)$. Pursuing our aim to ensure that the influence of the HOL delay has a dominant effect when the urgency of packet transmission is high and, vice versa, when the HOL delay is low, the weight $a_k(n)$ has an exponential influence on the utility. In order to calculate the utility value, each user has to feed back to the BS only the parameter $d_{\mathrm{HOL},k}(n)$, while the achievable transmission rate is calculated using CSI, gathered at the BS.

Due to differences in sensitivity to packet delays, the weights for delay-sensitive and for delay-tolerant traffic are adapted differently. Regardless of the traffic type, the actual QoS parameters of delay-sensitive users are always used, thus enabling the actual provision of best-effort service for delay-tolerant users, and preferential treatment of delay-sensitive users.

**Figure 4 The reduction of the search space for M = 4 and $\rho_{max} = 0.5$ in terms of (a) number of available users and (b) percentage of the full set of users**.

***Weight adaptation for delay-sensitive traffic*** It proves important to keep the transmission rate above the threshold level, and the packet end-to-end delay under the defined deadline for delay-sensitive applications (Figures 1b and 2). However, user satisfaction does not increase if we further decrease the delay, or increase the throughput. Therefore, the objective is to ensure that the delay is kept just under the threshold level and that the throughput is kept just above the threshold level, and hence, to optimize the utility while also preventing excessive use of resources for delay-sensitive applications.

For each delay-sensitive application, the minimum average throughput threshold $\overline{r_{min,k}}$ and the packet waiting time deadline $D_{max,\ k}$ are set according to the application characteristics. Note that the end-to-end delay consists of the time the packet spends in a queue (scheduling delay) and the time required for transmission across the network. Considering the variation in scheduling delay, the deadline has to be set proportionately lower than the difference between the required end-to-

end and transmission delays, in order to prevent the occasional deadline violation resulting in end-to-end delay violation. Therefore, the parameter $D_{max,\ k}$ does not present the absolute upper bound for the scheduling delay, yet only a reference point used for weight adaptation. Furthermore, as the transmission delay is a varying network-dependent value, the algorithm has to be able to support the adaptation of the waiting time deadline in order to constantly guarantee that the end-to-end delay requirements are met.

The weights are adapted periodically, based on the average QoS level, and calculated separately for scheduling delay and transmission rate. QoS level is calculated using the following equations:

$$\text{QoS}_{r,k} = \frac{\overline{r_k(n)}}{\overline{r_{min,k}}}, \tag{6}$$

$$\text{QoS}_{d,k} = \frac{D_{max,k}}{d_{\text{HOL},k}(n)}, \tag{7}$$

where $\overline{r_k(n)}$ and $\overline{d_{\mathrm{HOL},k}(n)}$ are calculated using an exponential moving average (EMA) function with forgetting factor $\alpha$, which defines the level of influence of the older values:

$$\overline{r_k(n)} = (1 - \alpha_r) \cdot \overline{r_k(n-1)} + \alpha_r \cdot r_k(n-1), \qquad (8)$$

$$\overline{d_{\mathrm{HOL},k}(n)} = (1 - \alpha_d) \cdot \overline{d_{\mathrm{HOL},k}(n-1)} + \alpha_d \cdot d_{\mathrm{HOL},k}(n-1). \quad (9)$$

Note that the average HOL delay is updated only if the user was selected in the previous frame. The values of parameters $\alpha_r$ and $\alpha_d$ are not equal–the scheduling algorithm exploits multiuser diversity. Therefore, the long-term average is more important for the transmission rate, which means that $\alpha_r$ should have a lower value. On the other hand, the delay has to be constantly kept under the deadline; hence, $\alpha_d$ should have a higher value.

While the individual average QoS level is used to provide the required QoS level, the fairness in resource allocation is provided with the use of a relative QoS level in relation to other users using the same traffic type. The intra-application user's QoS level is used to define the incrementation/decrementation step for the weight adaptation, and is calculated as the ratio of the user's individual QoS level to the averaged QoS level of all users using the same application type:

$$\mathrm{QoS}_{\mathrm{intra},k} = \frac{\mathrm{QoS}_k}{\frac{1}{K_{\mathrm{RT},i}} \cdot \sum_{k' \in K'} \mathrm{QoS}_{k'}}; \quad k \in K', \qquad (10)$$

where $K'$ is a subset of users using the same application type (e.g. the subset of VoIP users) and $K_{\mathrm{RT},\,i} = m(K')$, i.e. the cardinality of $K'$. The parameter $\mathrm{QoS}_{\mathrm{intra}}$ is calculated separately for the transmission rate and the HOL delay ($\mathrm{QoS}_{d\_intra}$ and $\mathrm{QoS}_{r\_intra}$).

Using these parameters, the weights for delay-sensitive (i.e. real-time (RT) users) are adapted as follows:

$$a_k(n) = \begin{cases} a_k(n-1) + \Delta a / \mathrm{QoS}_{d\_\mathrm{intra},k}; & \text{if} \quad \mathrm{QoS}_{d,k} < 1 - G_{\mathrm{RT}} \\ a_k(n-1) - \Delta a \cdot \mathrm{QoS}_{d\_\mathrm{intra},k}; & \text{if} \quad \mathrm{QoS}_{d,k} > 1 + G_{\mathrm{RT}} \\ a_k(n-1); & \text{otherwise} \end{cases}, \quad (11)$$

$$b_k(n) = \begin{cases} b_k(n-1) + \Delta b / \mathrm{QoS}_{r\_\mathrm{intra},k}; & \text{if} \quad \mathrm{QoS}_{r,k} < 1 - G_{\mathrm{RT}} \\ b_k(n-1) - \Delta b \cdot \mathrm{QoS}_{r\_\mathrm{intra},k}; & \text{if} \quad \mathrm{QoS}_{r,k} > 1 + G_{\mathrm{RT}} \\ b_k(n-1); & \text{otherwise} \end{cases}, \quad (12)$$

where $\Delta a$ and $\Delta b$ are positive constants defining the basic step for weight adaptation. The weights $a_k$ and $b_k$ are positive parameters initially set to value 1. The users recording lower satisfaction levels (i.e. lower intra-application QoS levels) are assigned a higher weight increment (or lower priority decrement), which results in better fairness properties of the algorithm. Note that the prerequisites $a_k(n) > 0$ and $b_k(n) > 0$ need to be always

fulfilled. The parameter $G_{\mathrm{RT}}$ is a guard interval, determining the responsiveness of the scheduling metric, and has the following range: $0 < G_{\mathrm{RT}} < 1$.

**Weight adaptation for delay-tolerant traffic** Due to the 'elastic' nature of the delay-tolerant BE traffic and its high adaptivity to delay and bandwidth (Figure 1a), the priority weights for such applications are adapted according to the average QoS level of the delay-sensitive users, instead of the individual QoS levels of BE users.

For BE applications, the intra-application QoS level is calculated only in terms of the transmission rate, given that this is the appropriate performance measure for such traffic:

$$\mathrm{QoS}_{\mathrm{BE},k} = \frac{\overline{r_k(n)}}{\frac{1}{K_{\mathrm{BE}}} \cdot \sum_{k'' \in K''} \overline{r_{k''}(n)}}; \quad k \in K''. \qquad (13)$$

$K''$ is the subset of BE users and $K_{\mathrm{BE}} = m(K'')$ is the cardinality of $K''$. As for the RT users, the intra-application of QoS level is used to define the incrementation/decrementation step for the adaptation of the weight $b_k$. The incrementation/decrementation step for the delay-dependent weight $a_k$ is constant and equals $\Delta a$:

$$a_k(n) = \begin{cases} a_k(n-1) + \Delta a; & \text{if} \quad \overline{\mathrm{QoS}_{d\_\mathrm{RT}}} > 1 + G_{\mathrm{BE}} \\ a_k(n-1) - \Delta a; & \text{if} \quad \overline{\mathrm{QoS}_{d\_\mathrm{RT}}} < 1 - G_{\mathrm{BE}} \\ a_k(n-1); & \text{otherwise} \end{cases}, \quad (14)$$

$$b_k(n) = \begin{cases} b_k(n-1) + \Delta b / \mathrm{QoS}_{\mathrm{BE},k}; & \text{if} \quad \overline{\mathrm{QoS}_{d\_\mathrm{RT}}} > 1 + G_{\mathrm{BE}} \\ b_k(n-1) - \Delta b \cdot \mathrm{QoS}_{\mathrm{BE},k}; & \text{if} \quad \overline{\mathrm{QoS}_{d\_\mathrm{RT}}} < 1 - G_{\mathrm{BE}} \\ b_k(n-1); & \text{otherwise} \end{cases}, \quad (15)$$

where $\overline{\mathrm{QoS}_{d\_\mathrm{RT}}}$ is the average value of parameters $\mathrm{QoS}_{d,k}$ from all RT users in the network:

$$\overline{\mathrm{QoS}_{d\_\mathrm{RT}}} = \frac{1}{K_{\mathrm{RT}}} \cdot \sum_{k=1}^{K_{\mathrm{RT}}} \mathrm{QoS}_{d,k}. \qquad (16)$$

A guard interval $G_{\mathrm{BE}}$ is also considered, although its value is not necessarily equal to $G_{\mathrm{RT}}$. The adopted approach allows an efficient allocation of available resources, since the priority of BE users is increased when, on average, RT users are experiencing high levels of QoS and decreased when available resources need to be assigned to RT users in order to provide the required level of QoS.

Optimal transmission scheme selection assuming zero-forcing receivers

Once the set of spatially multiplexed users is determined, the optimal transmission modes are selected for each user, using a recursive procedure at the BS that takes into account the user's estimated SNR after the signal detection, the properties of the available transmission modes, and the maximal BER requirements for

each traffic type. As the algorithm foresees the utilization of a linear ZF receiver, the SNR for the *i*th user after the detection, can be calculated analytically, as explained in [20–equations (1) to (7), 21].

Next, the approach proposed in [20] is adopted. If it proves impossible to meet the target BER constraint for all users sharing the same resources, we remove the user with the lowest utility in order to further decrease the CCI and hence, improve the conditions for the remaining users. This procedure is repeated until the required transmission reliability may be provided to all users, and then the optimal transmission mode is assigned to each user.

### Wireless system model and algorithm parameters

In our simulations, the base station is equipped with $M = 4$ antennas. The channel is assumed to be static for the duration of one frame and changes independently in the next frame. Perfect CSI at the BS is assumed. Channel coefficients for each user follow the Rayleigh distribution. As there are no recommendations for multiuser MIMO channel models, we defined a MIMO channel for each user and used the same distribution for all users in order to limit the impact of different channel characteristics on the performance evaluation of the proposed resource allocation scheme. A simplistic channel model is used in order to limit the effect of advanced channel model parameters, so the contribution of the scheduling metric to the system performance could be isolated. The effect of advanced propagation models, such as the COST 259 [22] and COST 273 [23] models, on the simulation results as well as the addition of Ricean distribution for channel coefficients of certain users and Kronecker correlation model, often used in MIMO systems, still have to be examined. However, it is expected that the performance of the proposed scheme, relative to the performance of the existing resource allocation schemes, will not change drastically, as this would affect each of them in the same manner. Furthermore, the importance of the proposed interference mitigation scheme would become even more significant in the system where users' channels would be more correlated.

Three different traffic types are taken into consideration, namely, VoIP, video streaming and BE traffic. Inside the cell with normalized radius $r = 1$, the users are located on $n$ equidistantly distributed virtual rings. Three users, each using a different traffic type (red circles depict VoIP users, green squares video streaming users and yellow diamonds depict BE users), are located on each ring (Figure 5); hence, $n = K/3$ and each traffic type is represented with $K/3$ users. The distance between the nearest ring and the BS is always $d = 0.1r$. Such a user distribution is chosen to eliminate the influence of non-uniform geographic distribution of applications inside the cell on the performance comparison of different resource allocation algorithms, which is the focus of this research.

We assume that all users transmit their data using the same normalized power $P°$, defined in such a manner that, in the interference-free channel, the edge-cell users can on average transmit their data using the most robust transmission mode available in the system. Using the proposed power control, we actually set the required average SNR at the edge of the cell. Nonetheless, the instantaneous SNR depends on the channel realization in each frame. Furthermore, the path loss exponent equals two. Applying different path loss exponent would only modify the SNR range inside the cell, or change the cell radius, if the SNR range was kept constant. Taking into account the assumed ring distribution and the path loss exponent, the difference in signal strength between the nearest and the furthest ring equals 20 dB.

The packets arrive in the queues at a constant rate $R_i$. The assumed arrival rates are; $R_{\mathrm{VoIP}} = 128$ kbits/s for VoIP traffic, $R_{\mathrm{VS}} = 384$ kbits/s for video streaming traffic and $R_{\mathrm{BE}} = 256$ kbits/s for BE traffic. The target BER values are $\mathrm{BER}_{\mathrm{RT\_max}} = 10^{-3}$ for RT traffic and $\mathrm{BER}_{\mathrm{BE\_max}} = 10^{-11}$ for BE traffic. For simulation purposes, we set the bandwidth to $B = 2$ MHz, while a time division duplex (TDD) system with frame duration $T_{\mathrm{f}} = 5$ ms is assumed. The ratio between the uplink and downlink shares in one time frame is taken from the IEEE 802.16-2005 communication standard [24], and is $T_{\mathrm{UL}}/T_{\mathrm{DL}} = 18/29$.

The set of available transmission modes is also taken from [24]. Nine transmission modes (QPSK, 16QAM and 64QAM modulations in combination with convolutional coding (CC) and a Reed-Solomon block encoder) are considered. The performance requirements for selected transmission modes in the AWGN channel, in terms of SNR thresholds for achieving the desired BER, are listed in Table 1. The results were obtained with Monte Carlo simulation.

### Performance analysis

The scheduling metric parameters used in simulations have the following values:

The packet waiting time deadline is set to $D_{\mathrm{max\_VoIP}} = 75$ ms for VoIP traffic and $D_{\mathrm{max\_VS}} = 150$ ms for video streaming traffic.

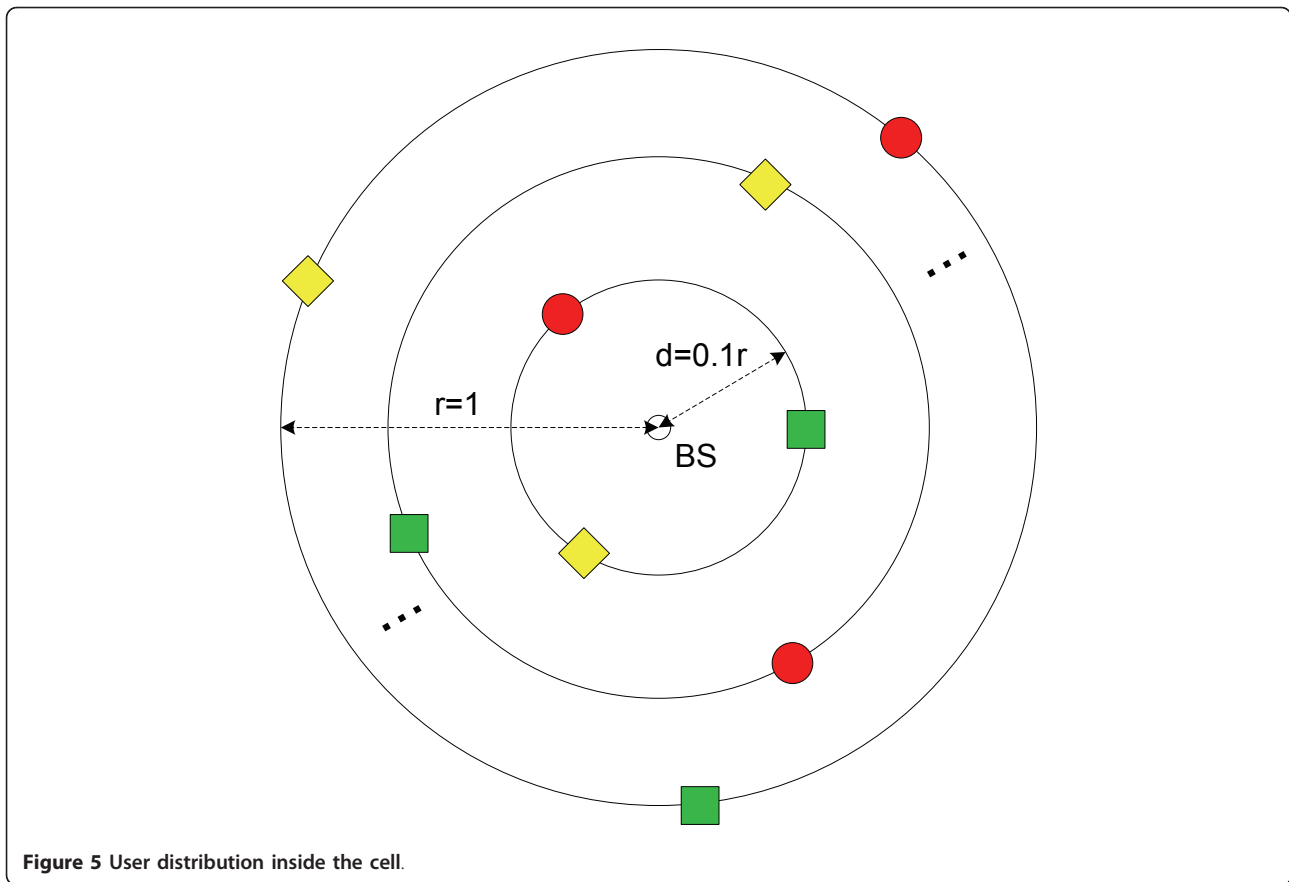The transmission rate threshold $\overline{r_{\mathrm{min},k}}$ is defined with the average arrival rate $R_k$.

Forgetting factors in EMA function are set to $\alpha_d = 0.6$ and $\alpha_r = 0.1$.

Basic weight adaptation step is set to $\Delta a = \Delta b = 0.02$.

Guard intervals are set to: $G_{\mathrm{RT}} = 0.2$ and $G_{\mathrm{BE}} = 0.1$.

Weights are adapted in every twentieth frame.

**Figure 5 User distribution inside the cell**.

Joint optimization of parameters $\alpha_d$, $\alpha_r$, $\Delta a$, $\Delta b$, $G_{RT}$ and $G_{BE}$ may be achieved with mathematical tools; however, the problem becomes very complex at a higher number of parameters. Therefore, we adopted a greedy approach, where parameters were tuned successively, based on test simulations.

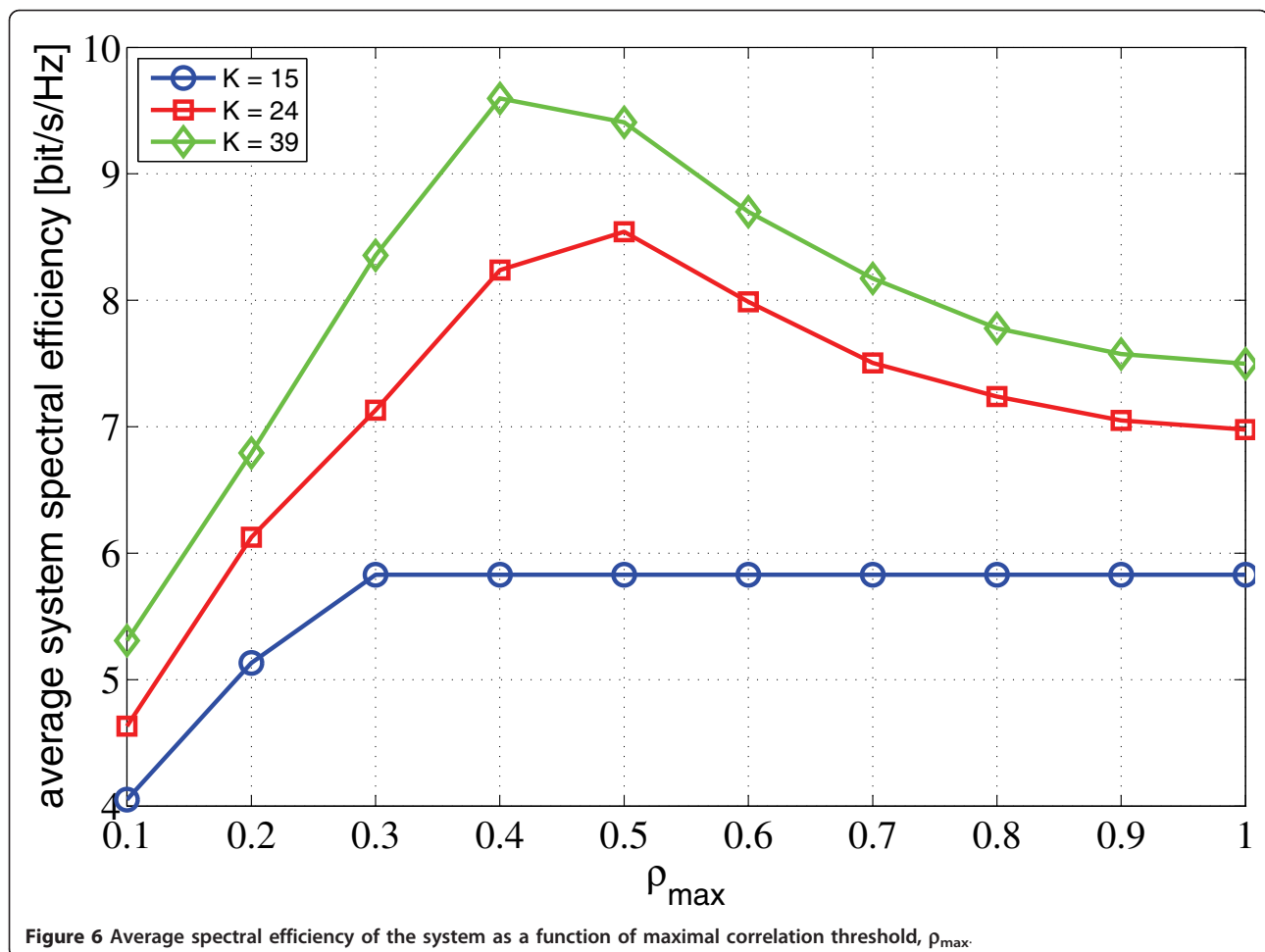Next, the optimal value for the maximal correlation threshold $\rho_{max}$, used in the proposed CCI mitigation technique, was investigated. The average spectral efficiency of the system as a function of $\rho_{max}$ depends on the number of users ($K$) in the system (Figure 6).

Simulations show that, for the system model assumed in the simulations, the optimal value is $\rho_{max} = 0.5$. Lower values of $\rho_{max}$ allow better CCI mitigation; however, it is more difficult to find the set of users not violating the maximal correlation condition, and therefore, fewer users are able to share common resources. In contrast, if $\rho_{max}$ is higher, signal distortion due to CCI becomes too high.

Due to low traffic load at $K = 15$, the selection of $\rho_{max}$ does not have an effect on the efficiency as long as $\rho_{max} \geq 0.3$, since the system is able to serve all the users efficiently, even under high CCI. With larger number of users in the system, the traffic load, as well as the multiuser diversity, becomes greater. Hence, it is easier to find the set of less correlated users. Consequently, an optimal value of correlation threshold $\rho_{max}$ can be determined. In theory (sufficient system capacity), the optimal value of $\rho_{max}$ would decrease continuously by increasing the number of users. However, in the assumed system, the traffic queues cannot be kept stable at $K = 39$, as will be seen later, therefore, the optimal value is $\rho_{max} = 0.5$ and this value will be used in further analysis.

**Table 1 Available transmission modes and performance requirements for AWGN channel in terms of SNR threshold [26 - Figure thirty-five]**

| Transmission mode | Spectral efficiency [bit/s/Hz] | SNR threshold [dB] (BER < 10⁻³) | SNR threshold [dB] (BER < 10⁻¹¹) |
|---|---|---|---|
| QPSK 1/2 | 0.937 | 2.65 | 4.15 |
| QPSK 2/3 | 1.250 | 4.40 | 5.85 |
| QPSK 3/4 | 1.406 | 5.30 | 6.60 |
| 16QAM 1/2 | 1.875 | 7.35 | 8.95 |
| 16QAM 2/3 | 2.500 | 10.10 | 11.55 |
| 16QAM 3/4 | 2.812 | 11.25 | 13.05 |
| 64QAM 2/3 | 3.749 | 14.70 | 16.40 |
| 64QAM 3/4 | 4.218 | 16.40 | 18.25 |
| 64QAM | 5.624 | 21.35 | 23.45 |

**Figure 6 Average spectral efficiency of the system as a function of maximal correlation threshold, $\rho_{max}$.**
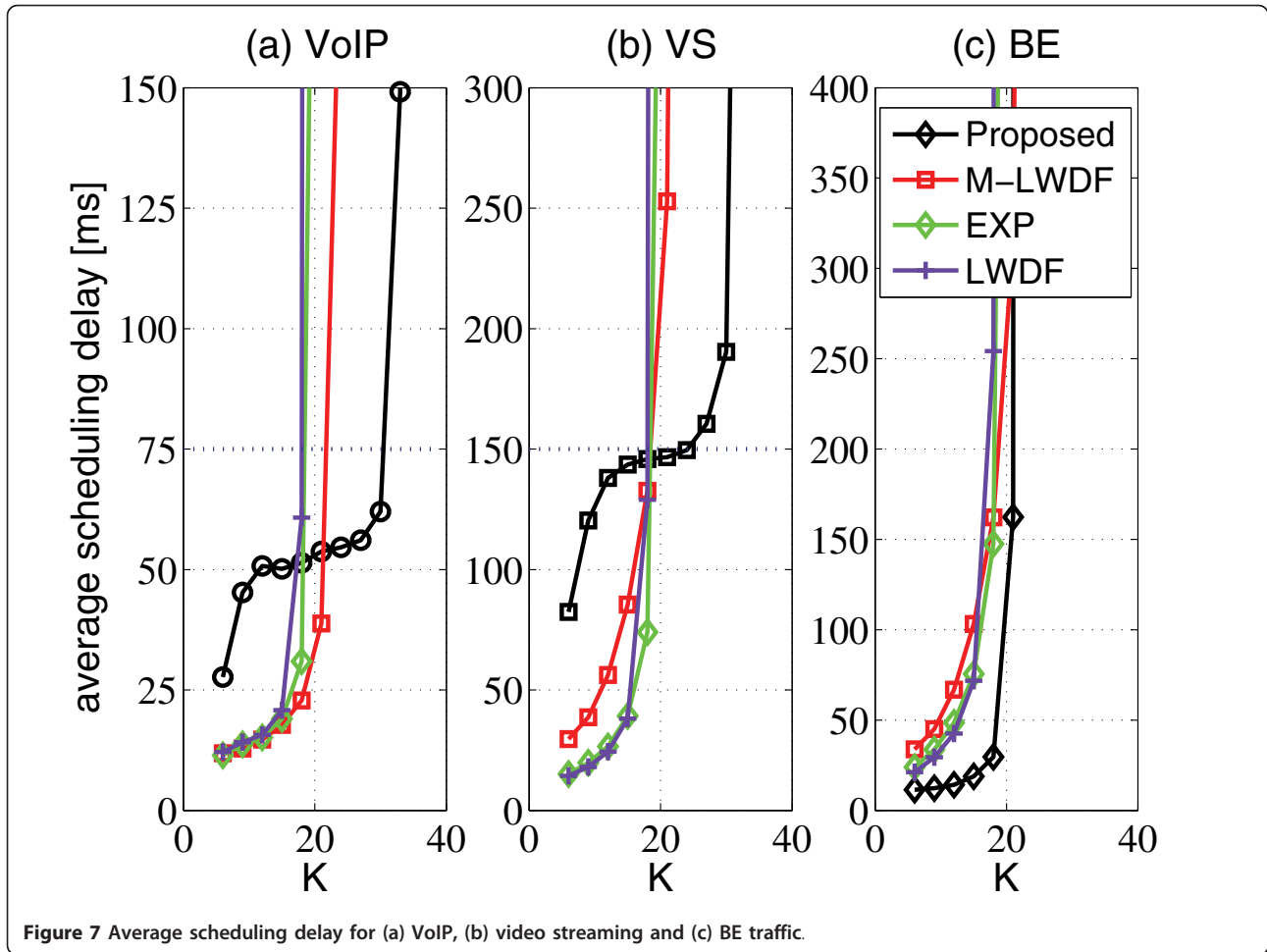
### Comparison of scheduling metrics

The efficiency of the proposed adaptive scheduling metric was evaluated by comparing the existing metrics, namely PFS, LWDF, M-LWDF and EXP. Figure 7 depicts the average scheduling delay for different traffic types. Note that the PFS algorithm performs very poorly, since it only considers the channel state and has no mechanism for QoS provision for delay-sensitive users. Although a certain fairness criterion is considered, the PFS rule often assigns resources to users with good channel conditions. Cell-edge users thus experience low service quality, and the average performance level deteriorates significantly since the results are averaged over the entire set of users using the same application. The average delay is too high and is thus not depicted in Figure 7. As expected, M-LWDF and EXP rules provide the best performance of all the existing scheduling metrics.

The simulations show that the use of proposed adaptive scheduling metric enables the queues of RT users to be kept stable for a higher number of active users than the use of other metrics. For VoIP users, the average scheduling delay is kept below the chosen deadline until $K > 30$, while for video streaming users the deadline is exceeded at $K > 24$, although it is kept at a reasonably low value at $K = 30$. Having in mind a particular level of adaptivity for such traffic (Figure 1c), we can say that a satisfactory level of QoS is achieved even at such a value of $K$. An additional consequence of the weight adaptation is the fact that, at low values of $K$, the average delay is closer to the deadline when the proposed metric is used, which is exactly what we sought to achieve. Adopting the proposed approach enables better utilization of radio resources, since more resources may be assigned to BE users, while maintaining the same QoS level for RT users. However, at high $K$, the adaptation of weights based on the QoS levels of RT users results in more significant deterioration of the QoS for BE users than is the case with other metrics. This can be seen clearly in Figure 8, which depicts the average user throughput for different traffic types. The upper bound of the average user throughput is defined as the average traffic arrival rate. Moreover, same conclusions can be extracted from both Figures 8 and 7; however, different performance measure is applied.

**Figure 7 Average scheduling delay for (a) VoIP, (b) video streaming and (c) BE traffic.**

In order to establish whether the proposed scheduling metric yields higher system level efficiency than the existing metrics, it is important to determine whether the performance increase for RT users prevails over the performance decrease for BE users at high values of $K$. This can be evaluated by comparing the average system utilities, which actually define the system level efficiency. For that purpose, we have used the approximation of utility curves as a function of transmission rate, as presented in Figure 1. The utility curve of elastic BE applications is described using the following exponential function:

$$u_{BE}(r) = 1 - e^{-\left(\frac{c}{\lambda_{BE}} \cdot r\right)}. \tag{17}$$

The units for the packet arrival rate $R_{BE}$ and transmission rate $r$ are kbits /s. The parameter $c$ is used to ensure that the utility approaches 1 at $r = R_{BE}$. A value of 4.61 is chosen for $R_{BE} = 256$ kbit /s.

The utility curves of RT applications are described using a sigmoid *arctan* function:

$$u_{RT}(r) = \frac{1}{\pi} \cdot \arctan(c_1 \cdot (r - c_2)) + 0.5. \tag{18}$$

The parameter $c_1$ defines the curve gradient, while the parameter $c_2$ defines the value of $r$ at which the utility reaches the value 0.5. Based on the packet arrival rates ($R_{VS} = 384$ kbit/s and $R_{VoIP} = 128$ kbit/s), the values of parameters $c_1$ and $c_2$ are: $c_1 = 0.3$ and $c_2 = 350$ kbit /s for video streaming traffic and $c_1 = 100$ and $c_2 = 127.5$ kbit /s for VoIP traffic, simulating a step function. The approximated utility curves are depicted in Figure 9.

The average utility, calculated using the approximated utility functions for individual traffic type and for the entire system is shown in Figures 10 and 11, respectively. The simulation results confirm that, despite a certain performance deterioration of BE traffic due to the increased stability of the queues of RT users, the average system utility is significantly higher if the proposed adaptive scheduling metric is applied. With the proposed approach, we efficiently exploit the elastic nature of BE traffic, since a satisfactory level of service can still be provided even if the throughput is decreased,

**Figure 8 Average user throughput for (a) VoIP, (b) video streaming and (c) BE traffic**.

whereas the utility of RT users in such a case decreases significantly. When the traffic load becomes too high and the system cannot keep the queues of RT users stable, even with the proposed weight adaptation, the aggregate utility, obtained with other metrics, is higher due to better conditions of BE users; however, this is just a theoretical result, since the call admission control (CAC) protocol will never accept such a number of connections due to the low QoS of RT users.

An additional objective of the proposed scheduling metric is to achieve good fairness properties. While, for homogeneous data traffic, fairness properties can be described well using different fairness indices, calculated on the basis of the fraction of resources each user is allocated, this is not the case for delay-sensitive RT applications. In such a case, the QoS level, expressed in terms of average utility, is a more appropriate measure. Figure 12, depicting the average utility for individual users at different number of users in the systems (each user is represented with a circle), shows that the same QoS level is provided to all VoIP and video streaming users, regardless of their location inside the cell. A higher level of deviation is evident for BE users. However, the average transmission rate assigned to BE users (Figure 13) shows that the deviation decreases.

The reason lies in the gradient of the utility curve at lower throughputs. The fairness levels for BE users were estimated, using the Jain fairness index [25] as a quantitative measure:

$$f(\bar{r}) = \frac{\left| \sum_{i=1}^{N} \bar{r}_i \right|^2}{N \sum_{i=1}^{N} \bar{r}_i^2}, \tag{19}$$

where $N$ is a population of BE users, i.e. $N = K_{BE} = K/3$. The results for different values of $K_{BE}$ are summarized in Table 2. Very good fairness properties are achieved for BE users as long as the traffic load is not too high (i.e. until $K_{BE} > 10$ or as long as RT queues are kept stable). When the traffic load becomes too high, the fairness index decreases; however, as explained hereinabove, such a scenario cannot occur in the given system due to the CAC protocol.
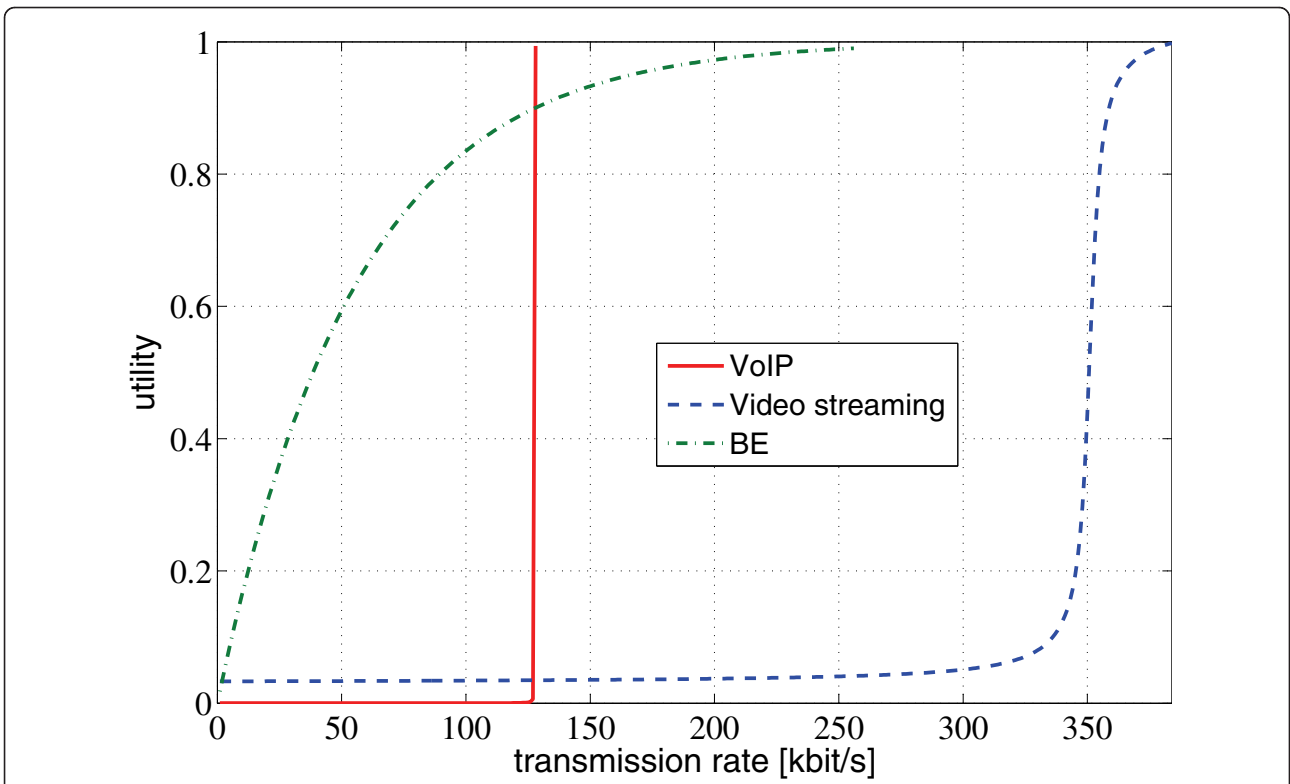
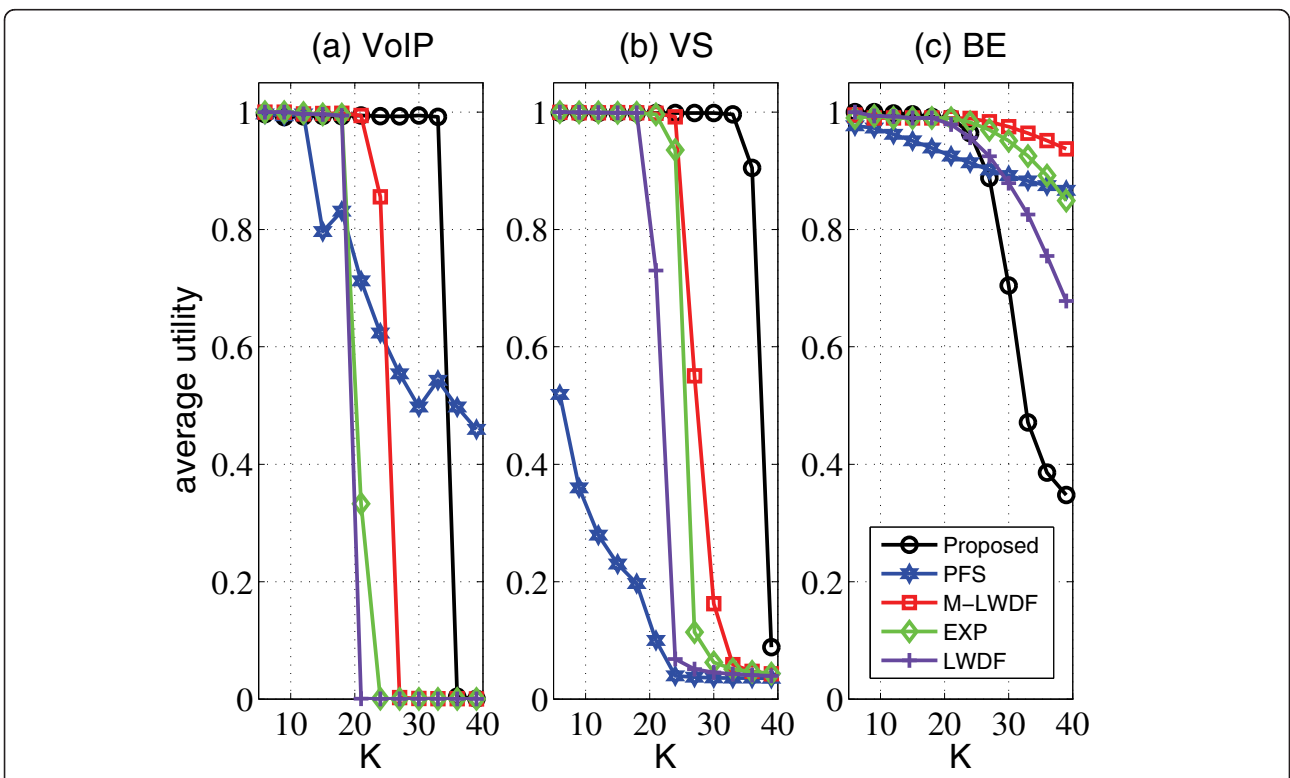**Figure 9 Approximated utility curves for different traffic types**.



**Figure 10 Average utility for (a) VoIP, (b) video streaming and (c) BE traffic**.
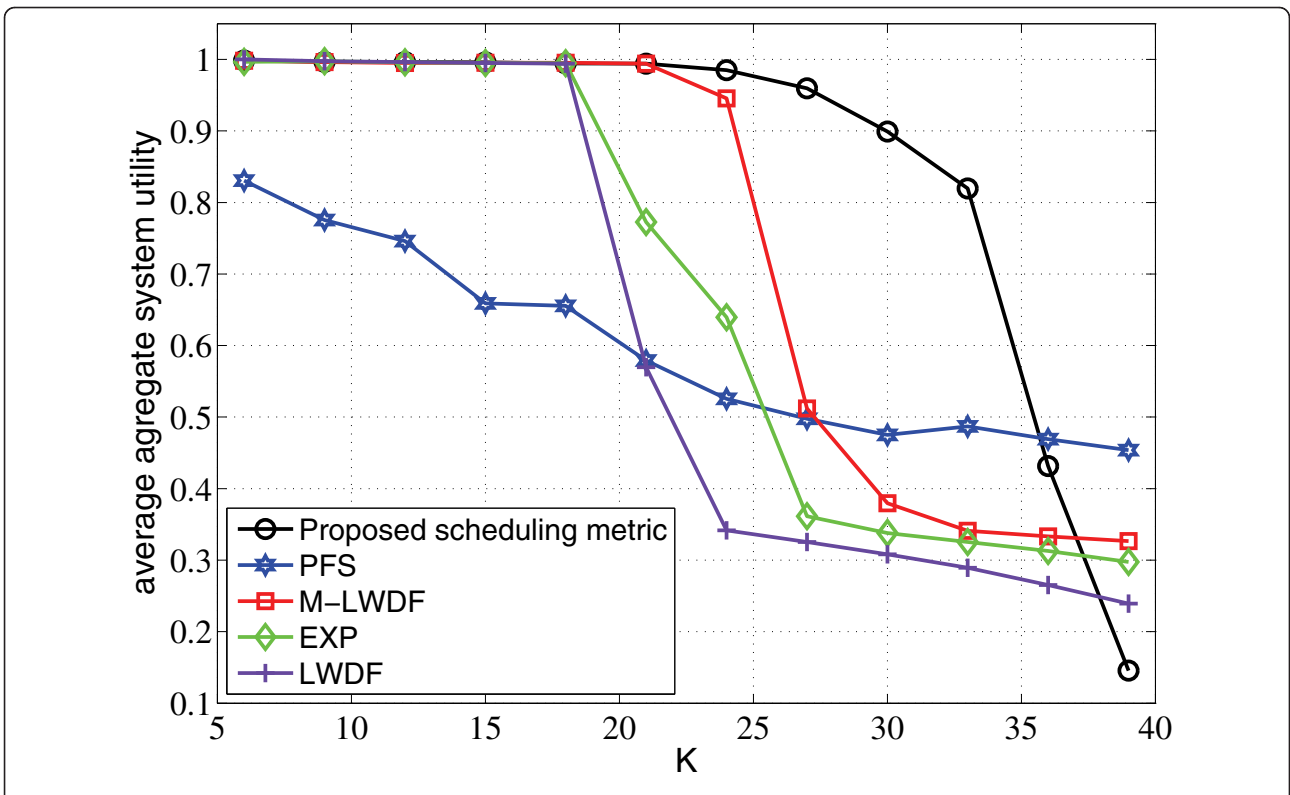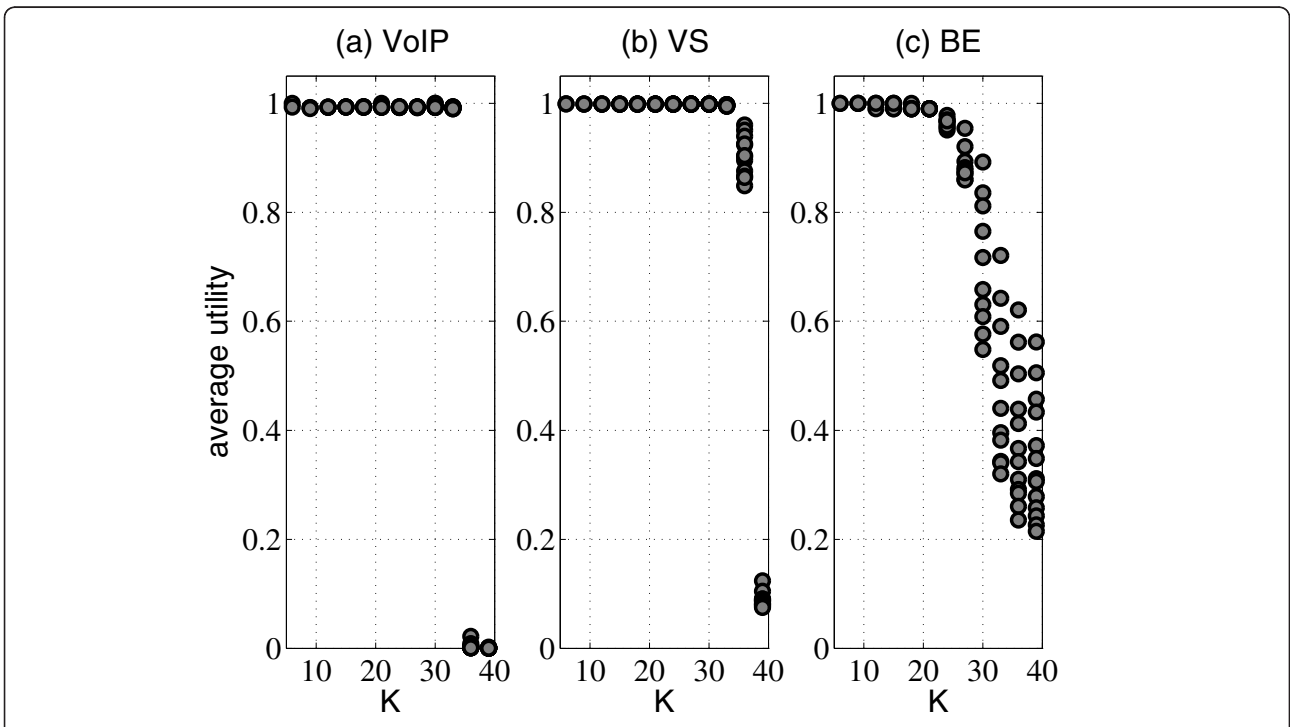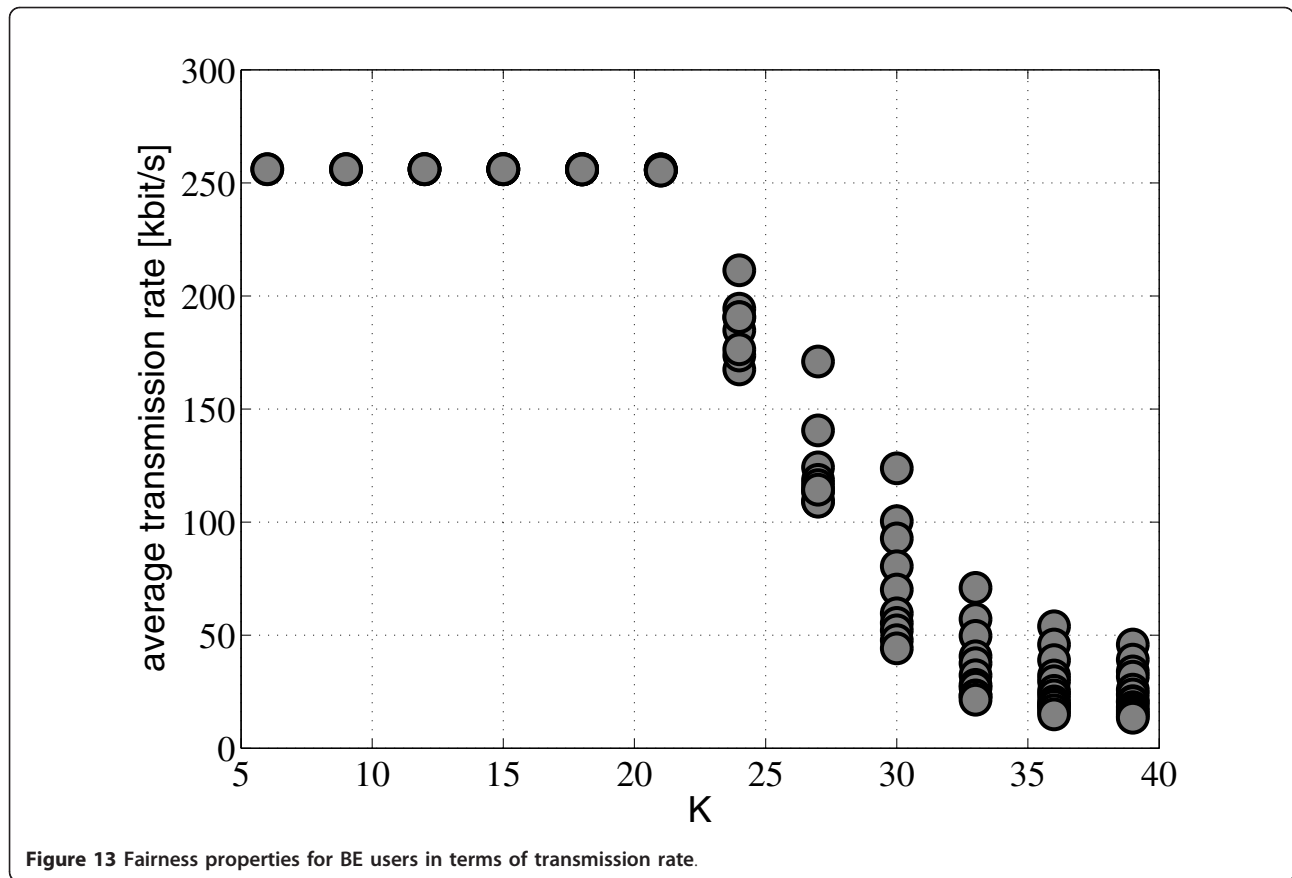
**Figure 11 Average aggregate system utility.**



**Figure 12 Fairness properties of the proposed scheduling metric in terms of average utility for (a) VoIP, (b) video streaming and (c) BE users.**

**Figure 13 Fairness properties for BE users in terms of transmission rate**.

### Virtual MIMO vs. SISO comparison

Finally, the actual benefit of using multiple antennas at the BS is presented, allowing for spatial multiplexing of users, in contrast to allocating resources to a single user, as it is implemented in conventional single antenna systems. Figure 14 depicts the average system spectral efficiency, achieved in a case $M = 1$ (SISO system) and in a case $M = 4$ (virtual MIMO system). It can be seen clearly that an increase in the linear capacity for the factor 4 is actually achieved.

### Conclusions

A scheduling algorithm for multiuser MIMO uplinks, enabling spatial multiplexing of users to be supported, is presented. The objective of the proposed algorithm has been to optimize the resource allocation in heterogeneous systems with diverse QoS requirements while, at the same time, providing fair resource allocation. The CCI mitigation, arising from the correlation of spatially

multiplexed users sharing common resources, is achieved efficiently with the use of a parameter $\rho_{\max}$, which defines a maximal correlation between any pair of spatially multiplexed users. The complexity of incremental user selection also decreases with the adopted approach due to search space reduction.

The main contribution of this work is the design of an adaptive channel and queue-aware, utility-based scheduling metric, the advantage of which lies in the periodic adaptation of priority weights based on the application of specific characteristics. Compared with the existing utility-based scheduling metrics, the results show a considerable performance improvement in terms of aggregate system utility, especially under higher traffic loads. The proposed adaptation is especially beneficial for RT users, since it allows excellent control over their QoS parameters. Benefits for BE users are observed at a lower number of users in the system whereas, at a higher number of users, their QoS level deteriorates at the expense of performance

**Table 2 Jain index of fairness for the proposed scheduling metric**

| $K_{BE}$ | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | $\overline{f(\bar{r})}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $f(\bar{r})$ | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.99 | 0.98 | 0.90 | 0.86 | 0.86 | 0.86 | 0.953 |

**Figure 14 Average system spectral efficiency for M = 1 and M = 4.**

improvement for RT users. Nevertheless, the elastic nature of BE traffic still enables a satisfactory QoS level.

In addition, the simulation results confirm that the proposed scheduling metric yields very good fairness properties in terms of user QoS levels, since users experience the same QoS levels at the cell-edge as users in the centre of the cell.

As the focus of the article is the comparison of the proposed metric with different existing utility-based scheduling metrics, several simplistic assumptions regarding propagation characteristics, channel model, user distribution or path loss exponent were made. Such assumptions were made in order to obtain a more straightforward comparison of different scheduling metrics, since we eliminate the influence of system-specific parameters that might distort the performance of scheduling algorithm efficiency.

### List of Abbreviations

BE: best-effort; BS: base station; CCI: co-channel interference; DPC: dirty paper coding; EXP: exponential; HOL: head-of-line; LWDF: largest weighted delay first; M-LWDF: modified LWDF; MIMO: multiple input multiple output; QoS: quality of service; PFS: proportionally fair scheduling; RT: real-time; SDMA: spatial domain multiple access; TDD: time division duplex; ZF: zero-forcing.

### Author details

[1]The Centre of Excellence for Biosensors, Instrumentation and Process Control - COBIK, Velika pot 22, SI-5250 Solkan, Slovenia [2]Department of Communication Systems, Jozef Stefan Institute, Jamova cesta 39, SI-1000 Ljubljana, Slovenia

### References

1. SG Glisic, *Advanced Wireless Communications* (John Wiley & Sons Inc., Hoboken, NJ, USA, 2007)
2. E Biglieri, R Calderbank, A Constantinides, A Goldsmith, A Paulraj, HV Poor, *MIMO Wireless Communications* (Cambridge University Press, Cambridge, UK, 2007)
3. C Fortuna, M Mohorcic, Trends in the development of communication networks: cognitive networks. Comput Netw. **53**(9), 1354–1376 (2009)
4. GJ Foschini, Layered space-time architecture for wireless communication in a fading environment when using multiple antennas. Bell Labs Tech J Autumn, 41–59 (1996)
5. D Gesbert, M Shafi, D Shiu, PJ Smith, A Naguib, From theory to practice: an overview of MIMO space-time coded wireless systems. IEEE J Sel Areas Commun. **21**, 281–302 (2003)
6. M Costa, Writing on dirty paper. IEEE Trans Inform Theory, **29**(3), 439–441 (1983)
7. T Yoo, A Goldsmith, On the optimality of multiantenna broadcast scheduling using zero-forcing beamforming. IEEE J Sel Areas Commun. **24**(3), 528–541 (2006)
8. G Dimić, ND Sidiropoulos, On downlink beamforming with greedy user selection: performance analysis and a simple new algorithm. IEEE Trans Sig Process. **53**(10), 3857–3868 (2005)
9. Z Shen, JG Andrews, RW Heath Jr, B Evans, Low complexity user selection algorithms for multiuser MIMO systems with block diagonalization. IEEE Trans Sig Process. **54**(9), 3658–3663 (2006)
10. M Shariat, AU Quddus, SA Ghorashi, R Tafazolli, Scheduling as an important cross layer operation for emerging broadband wireless systems. IEEE Comm Surveys Tutr. **11**(2), 74–86 (2009)
11. P Visawanath, DNC Tse, R Laroia, Opportunistic beamforming using dumb antennas. IEEE Trans Info Theory, **48**(6), 1277–1294 (2002)
12. R Knopp, P Humblet, Information capacity and power control in single cell multi-user communications, in *Proceedings of IEEE International Confeernce on Communication (ICC)*, vol. 1. (Seattle, WA, USA, 1995), pp. 331–335
13. A Stoylar, K Ramanan, Largest weighted delay first scheduling: large deviations and optimality. Ann Appl Prob. **11**(1), 1–48 (2001)

14. M Andrews, K Kumaran, K Ramanan, A Stoylar, P Whiting, Providing quality of service over a shared wireless link. IEEE Comm Mag. **39**(2), 150–154 (2001)

15. S Shakkottai, A Stoylar, Scheduling for multiple flows sharing a time-varying channel: the exponential rule. Bell Laboratories Technical Report, 2000

16. S Shenker, Fundamental design issues for the future internet. IEEE J Sel Areas Commun. **13**(7), 1176–1188 (1995)

17. P Liu, R Berry, ML Honig, Delay-sensitive packet scheduling in wireless networks, in *IEEE Wireless Communication and Networking 2003 (WCNC 2003)*, vol. 3. (New Orleans, LA, USA, 2003), pp. 1627–1632

18. YJ Zhang, KB Letaief, An efficient resource-allocation scheme for spatial multiuser access in MIMO/OFDM systems. IEEE Trans Commun. **53**(1), 107–116 (2005)

19. SY Tang, D Chieng, YC Chang, Uplink traffic scheduling with QoS support in broadband wireless access networks, in *Proceedings - MICC 2009: 2009 IEEE 9th Malaysia International Conference on Communication*, 623–628 (2009)

20. S Plevel, T Javornik, G Kandus, A recursive link adaptation algorithm for MIMO systems. AEU - Int J Electron Commun. **59**(1), 52–54 (2005)

21. RW Heath, S Sandhu, A Paulraj, Antenna selection for spatial multiplexing systems with linear receivers. IEEE Commun Lett. **5**(4), 142–144 (2001)

22. LM Correia, (ed), *Wireless Flexible Personalised Communications (COST 259 Final Report)* (John Wiley, Chichester, UK, 2001)

23. LM Correia, (ed), *Mobile Broadband Multimedia Networks, Techniques, Models and Tools for 4G* (Academic Press, Elsevier Ltd, UK, 2006)

24. IEEE Standard 802.16e-2005, Amendment to IEEE Standard for Local and Metropolitan Area Networks - Part 16: Air Interface for Fixed Broadband Wireless Access Systems - Physical and Medium Access Control Layers for Combined Fixed and Mobile Operation in Licensed Bands (2005)

25. RK Jain, DMW Chiu, WR Hawe, A quantitative measure of fairness and discrimination for resource allocation in shared computer systems. DEC-TR-301, Digital Equipment Corporation, Tech Rep (1984)

26. S Plevel, Adaptive multiple input multiple output wireless communication systems. PhD Thesis (in Slovene with english abstract), Ljubljana, SI, 2007