*Research Article*

# Throughput-Optimal Scheduling with Low Average Delay for Cellular Broadcast Systems

## Chan Zhou and Gerhard Wunder

*German-Sino Lab for Mobile Communications, Fraunhofer Institute for Telecommunications, Heinrich-Hertz-Institut, Einsteinufer 37, 10587 Berlin, Germany*

Correspondence should be addressed to Chan Zhou, chan.zhou@hhi.fhg.de

While a number of scheduling policies achieve the maximum throughput region, the average delay minimization problem for cellular broadcast systems still awaits its complete solution. To this end, we introduce a scheduling policy which decomposes the cross-layer delay optimization problem into two subproblems: allocation of physical resources and user priority management. The first subproblem is translated into a weighted sum rate maximization problem that can be efficiently solved for different channel models. The solution of the second subproblem determines the weight factors in the maximization problem expressing the priorities of users. For the latter subproblem we present a so-called *idle state prediction* algorithm minimizing our relevant delay measure. Analytical and simulative tools are used to show that the introduced scheduling policy provides both optimal throughput and low delay performance.

## 1. Introduction

Allocating limited resources at medium access control (MAC) layer and physical (PHY) layer among users (from now on briefly *scheduling*) is a fundamental problem in the design of next generation wireless systems. In general, a scheduling problem can be formulated as some kind of optimization problem, where the objective is to maximize/minimize some system performance measure under PHY layer constraints as well as quality of service (QoS) constraints on the MAC layer. One important performance measure is the *total system throughput* and it is therefore often considered as an objective of the optimization problem. *Throughput-optimal* scheduling policies are the policies, which can support any vector of arrival rates inside the ergodic achievable rate region [1, 2]. There exist quite a few scheduling policies which achieve this figure of merit [1, 3–9].

An important observation is that even though two different scheduling policies have the same throughput performance, they might significantly differ in, for example, their packet delay performance. Hence, in a system with random packet arrivals stored temporarily in queues, an enhanced performance criterion is to keep the queue lengths as short as possible so that the average packet delay of each user is minimized. One widely applied scheduling policy is the maximum weight matching scheduling (MWMS) policy which maximizes the sum of the rates weighted by the packet queue length of each user [3, 4, 10–12]. It was shown in [12] that the MWMS policy is delay-optimal for multiple-access channels. However, this result is based on the polymatroidal structure of the capacity region of multiple-access channels. For broadcast channels (BCs), MWMS is not delay-optimal even with symmetry assumptions. Motivated by this fact, Seong et al. introduced in [9] another throughput-optimal scheduling policy called queue proportional scheduling (QPS) which provides superior delay and fairness properties for the BC compared to MWMS. It minimizes the maximum draining time of the queueing system without new arrival. Based on the QPS policy, the delay region for such queueing systems can be characterized if the channel state is quasistatic [13].

A disadvantage of the approach in [9] is that the cost function is not directly related to average packet delay which

by Little's law can be calculated as the average queue length divided by the average arrival rate. Based on this expression, a scheduling problem can be formulated as a cross-layer optimization problem containing system parameters and constraints in PHY layer and MAC layer. A direct solution of such a problem involves a large number of optimization variables and is intractable. The focus of this paper is to provide a new approach to this problem.

*Contributions.* We first analyze some general characteristics of throughput-optimal scheduling policies. It is shown that they can be generally formulated as weighted sum rate maximization problems differing only in the choice of the weight factors. Furthermore, we prove that for throughput-optimal policies the weight factors are independent of the current channel state. Hence, the cross-layer scheduling problem is decomposed into two separate optimization subproblems:

(1) finding the optimal weight factors according to the queue states;

(2) solving the weighted sum rate maximization problem with respect to the instantaneous channel states.

Both subproblems are only coupled by the weight factors in the maximization problem. Interestingly, it was pointed out in [14, 15] that the solutions of such an optimization problem itself exhibit a layered structure with only limited degree of cross-layer coupling. Under some mild conditions, the complete optimization problem can be decomposed into several subproblems and the interfaces among them are quantified as the optimization variables coordinating the subproblems.

For Step 1 we introduce an iterative algorithm called *idle state prediction* (ISP) algorithm to obtain the optimal weight factors. This algorithm calculates the delay-optimal weight factors under the assumption that no new arrivals occur in the future by using the ergodic achievable rate region and the current queue state. Obviously, since we assume a dynamic scenario with random arrivals, the weight factors have to be recalculated in each time slot according to the updated queue state. Once the weight factors are fixed, the actual resource allocation is determined by maximizing the weighted sum of rates according to the instantaneous channel state. Note that we do not further comment on the weighted sum rate maximization problem in Step 2 for which there exist already efficient algorithms (multiple input multiple output (MIMO) [16–18], orthogonal frequency division multiplex (OFDM) [19, 20]). Simulations show that the delay performance can be significantly improved by the introduced scheduling policy.

The rest of this paper is organized as follows. In Section 2 we describe the system model and define the stability and delay measurement used in this paper. The characteristics of the throughput-optimal scheduling are analyzed in Section 3. Based on these characteristics, we introduce our parameter separation concept for the design of throughput-optimal policies. In Section 4 we present our scheduling policy for both static and dynamic channels.

The scheduler is evaluated through simulations in Section 5. Finally, we conclude with Section 6.

*Notations.* We use boldface letters to denote vectors and normal fonts with subscript are the elements of the vectors. $\|\mathbf{x}\|$ denotes the $l_1$-norm of the vector $\mathbf{x}$. The inequality between two vectors $\mathbf{x} \leq \mathbf{y}$ stands for $\mathbf{x}$ being componentwise smaller than or equal to $\mathbf{y}$. Furthermore we use $\lceil x \rceil$ to denote the smallest integer larger than $x$ and $\mathcal{A}^c$ is the complement of a set $\mathcal{A}$.

## 2. System Model

*2.1. PHY Layer.* We consider a single-cell downlink system in which a base station simultaneously supplies $M$ mobile users. The channel between the base station and each user is assumed to be constant within a time slot and varies from one time slot to another in an i.i.d. manner. The channel state of user $m$ in the $n$th time slot is denoted as $h_m(n) \in \mathcal{S}$, where $\mathcal{S}$ is an arbitrary countable or uncountable set, and all channel states of the user set $\mathcal{M} := \{1, \dots, M\}$ are collected in the vector $\mathbf{h}(n) \in \mathcal{S}^M$. Here, the set $\mathcal{S}$ is used to indicate that the general approach is not restricted to a specific transmission scheme. For example, in an MIMO system the channel state can be described as a matrix of complex channel gains such that $h_m(n) \in \mathbb{C}^{n_r n_t}$ where $n_r$, $n_t$ are the number of transmit and receive antennas at the base station and mobiles, respectively. Likewise, for an OFDM system the channel state can be defined as a vector of complex channel gains on each subcarrier $h_m(n) \in \mathbb{C}^K$, where $K$ is the number of subcarriers. In the $n$th time slot the data is transmitted through the channel at rate $\mathbf{r}(n) \in \mathbb{R}_+^M$ lying in the achievable rate region denoted as $\mathcal{C}(\mathbf{h}(n), \overline{P})$ with given sum power budget $\overline{P}$. For technical reasons we assume that the transmit rates $r_m(n)$ are uniformly bounded by some real constant $c_r > 0$.

Note that it is not relevant for the purposes of this paper in what way the achievable rate region is parameterized. It is just assumed that we are able to solve the following maximization problem:

$$\mathbf{r}(\boldsymbol{\mu}, \mathbf{h}(n)) := \underset{\widetilde{\mathbf{r}} \in \mathcal{C}(\mathbf{h}(n), \overline{P}),}{\arg \max} \boldsymbol{\mu}^T \widetilde{\mathbf{r}}, \qquad (1)$$

where $\boldsymbol{\mu} \in \mathbb{R}_+^M$ is the set of weight factors. The solution of (1) is a point on the convex hull of the achievable rate region $\mathcal{C}(\mathbf{h}(n), \overline{P})$. Observe that $\boldsymbol{\mu}$ also represents the normal vector of the convex hull at the point $\mathbf{r}(\boldsymbol{\mu}, \mathbf{h}(n))$ (see Figure 1). Then, the ergodic achievable rate region is given by

$$\mathcal{C}_{\text{erg}}(\overline{P}) := \bigcap_{\|\boldsymbol{\mu}\|=1} \{\widetilde{r}_1, \dots, \widetilde{r}_M : \boldsymbol{\mu}^T \widetilde{\mathbf{r}} \leq \boldsymbol{\mu}^T \mathbb{E}\{\mathbf{r}(\boldsymbol{\mu}, \mathbf{h}(n))\}\}, \quad (2)$$

which is a convex set [19, 20]. In this paper, we call the achievable rate region in time slot $n$ the *instantaneous achievable rate region* or just achievable rate region which is dependent on the current channel state $\mathbf{h}(n)$ and power constraint $\overline{P}$. The term *ergodic achievable rate region* is used for the rate region defined in (2) which is averaged over all channel states.
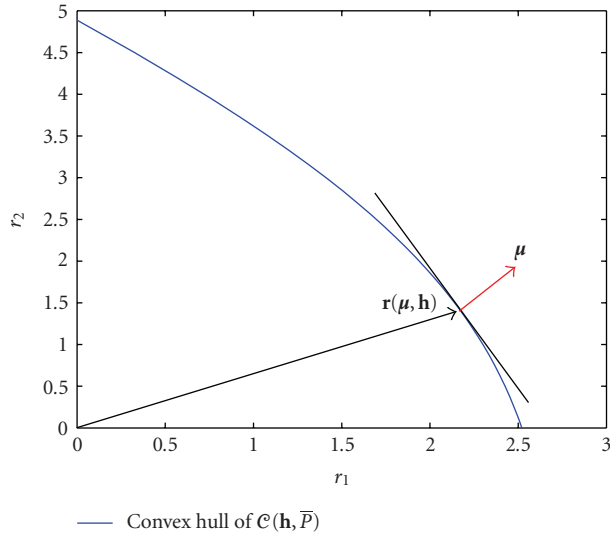
FIGURE 1: Ergodic rate region and expected rate allocation in a 2-user scenario. The solution of the weight maximization problem (1) is a point on the convex hull of the achievable rate region $\mathcal{C}(\mathbf{h}(n), \overline{P})$. The vector of weight factors $\boldsymbol{\mu}$ can be interpreted as the normal vector of the convex hull at the obtained point.

In order to show the applicability of our approach let us provide an example. We denote by $\mathcal{C}_{\text{OFDMA}}(\mathbf{h}(n), \overline{P})$ the achievable region of an orthogonal frequency division multiple access (OFDMA) system, where each subcarrier is exclusively assigned to one user. Due to the limited number of coding and modulation schemes only certain rates are achievable. Furthermore, there is a fixed power budget on each subcarrier $k$ denoted by $p_k$ and $\sum_k p_k \leq \overline{P}$. The achievable OFDMA region is defined as

$$\mathcal{C}_{\text{OFDMA}}(\mathbf{h}(n), \overline{P})$$
$$\equiv \bigcup_{\substack{\sum_{m=1}^{M} \theta_{m,k}=1, \forall k \\ \theta_{m,k} \in \{0,1\}}} \left\{ \mathbf{r} : r_m = \sum_{k=1}^{K} r'_{m,k}(\hat{h}_{m,k}, p_k) \theta_{m,k} \right\}, \quad (3)$$

where $\theta_{m,k} \in \{0,1\}$ is the indicator if user $m$ is mapped onto subcarrier $k$ and $r'_{m,k}(\hat{h}_{m,k}, p_k)$ is the rate of user $m$ on subcarrier $k$ with transmit power $p_k$ in one time slot (e.g., 2 bits for QPSK). The parameter $\hat{h}_{m,k}$ is the reported (quantized) channel state of user $m$ on subcarrier $k$. Due to these practical constraints $\mathcal{C}_{\text{OFDMA}}(\mathbf{h}(n), \overline{P})$ is a set of discrete rate points, nevertheless the solution in (1) achieves the points on the convex hull of $\mathcal{C}_{\text{OFDMA}}(\mathbf{h}(n), \overline{P})$ and the formulation of the ergodic achievable region $\mathcal{C}_{\text{erg}}(\overline{P})$ is still valid. A detailed description of this region can be found in [20].

### 2.2. MAC Layer.
Assuming that the transmission is time-slotted, data packets arrive randomly at the MAC and queue up in a buffer reserved for each user $m \in \mathcal{M}$. Simultaneously the data is read out from the buffers according to the system state, that is, the random channel state and the current queue lengths. Thus, the system can be modeled as a queueing

system with random processes reflecting the arrival and the departure of data packets.

Denoting the queue state of the $m$th buffer in time slot $n \in \mathbb{N}$ by $q_m(n)$ and arranging all queue states in the vector $\mathbf{q}(n) \in \mathbb{R}_+^M$ the evolution of the queueing system can be written as

$$\mathbf{q}(n + 1) = [\mathbf{q}(n) - \mathbf{r}(n) + \mathbf{a}(n)]^+, \quad (4)$$

where $[x]_m^+ = \max\{0, x_m\}$, for all $m \in \mathcal{M}$. $\mathbf{a}(n) \in \mathbb{R}_+^M$ is a random vector denoting the amount of arrival packets during the $n$th time slot and the vector $\mathbf{r}(n) \in \mathbb{R}_+^M$ is the amount of transmitted data. Without loss of generality we set the length of a time slot $T = 1$, so that $\mathbf{a}(n)$ and $\mathbf{r}(n)$ are equal to the arrival and transmit rate during the time slot $n$. We assume that the size of a data packet is constant and the sequence of arrival packets forms an i.i.d. sequence of variables over time. To simplify the notation we set the packet size to 1 without loss of generality. Further we make the technical assumption that the maximum numbers of arrival packets within one time slot are uniformly bounded by real constant $c_a > 0$. The vector of mean arrival rates is denoted as $\boldsymbol{\rho} = \mathbb{E}\{\mathbf{a}(n)\} \in \mathbb{R}_+^M$.

The transmit rate $\mathbf{r}(n)$ is determined by the applied scheduling policy. In this paper, we consider only stationary scheduling policies which are the mappings

$$\mathcal{P} : \mathcal{S}^M \times \mathbb{R}_+^M \longmapsto \mathbb{R}_+^M, \qquad (\mathbf{h}(n), \mathbf{q}(n)) \curvearrowright \mathbf{r}(n). \quad (5)$$

We further assume that the scheduling policies are only dependent on the proportion of the individual queue lengths and not dependent on the norm $\|\mathbf{q}(n)\|$; this covers most existing policies. The rate allocation according to the scheduling policy $\mathcal{P}$ is denoted as $\mathbf{r}^{\mathcal{P}}(\mathbf{h}(n), \mathbf{q}(n))$. Since both arrival rates $\mathbf{a}(n)$ and channel state $\mathbf{h}(n)$ are i.i.d., the evolution of the queueing system can be modeled as a discrete-time Markov chain with general state space.

### 2.3. A Cross-Layer Performance Measure.
The throughput region is defined as the set of all arrival rate vectors for which the Markov chain in (4) is stable in some sense [2]. There exist several relevant stability measures for Markov chains in literature, for example, strongly stable, weakly stable, recurrent. In this paper we resort to the definition of weak stability as in [21].

*Definition 1.* If for every $\epsilon > 0$ there is $B > 0$ and $N_0(\epsilon)$ such that for all $n > N_0(\epsilon)$, it follows $\Pr(\|\mathbf{q}(n)\| > B) < \epsilon$; then the Markov chain is *weakly stable*.

In contrast the definition of an unstable Markov chain is more subtle.

*Definition 2.* A Markov chain is said to be *uniformly transient* if there is a countable cover of $\mathbb{R}_+^M$ with uniformly transient sets, that is, there is $C < \infty$ with $\mathbb{E}\{\eta_\mathcal{A}\} \leq C$ for all $x \in \mathcal{A}$, where $\eta_\mathcal{A} = \sum_{n=1}^{\infty} \mathbb{I}(\mathbf{q}(n) \in \mathcal{A})$.

Let us introduce the following terminology: we call a vector of mean arrival rates $\boldsymbol{\rho}$ *stabilizable* (not stabilizable)

under a specific scheduling policy $\mathcal{P}$ when the corresponding queueing system driven by $\mathcal{P}$ is weakly stable (uniformly transient). It is well-known that any vector of arrival rates inside the ergodic achievable rate region is stabilizable and all other vectors of arrival rates are not stabilizable [1, 2]. Thus, a scheduling policy is called *throughput-optimal* if it keeps the system weakly stable for any vector of arrival rates $\boldsymbol{\rho}$ which lies in the ergodic achievable rate region.

Having defined throughput-optimal scheduling policies we now introduce a relevant cross-layer performance measure for average packet queueing delay. Consider the following quantity:

$$\overline{D}(N) = \frac{1}{M}\sum_{m=1}^{M} D_m(N) = \frac{1}{MN}\sum_{n=1}^{N}\sum_{m=1}^{M}\alpha_m q_m(n), \quad (6)$$

where the natural number $N \geq 1$ is the length of the observation time window and $\alpha_1,\dots,\alpha_m$ are positive real factors. If the factors are chosen such that $\alpha_m := 1/\rho_m = 1/\mathbb{E}\{a_m(n)\}$, and the limit $\lim_{N\to+\infty}(1/N)\sum_{n=1}^{N}\alpha_m q_m(n)$ exists and is equal to its stationary value, then $D_m(N)$ represents the average queueing delay of each packet of user $m$ as $N \to +\infty$ [22]. Note that even if the average arrival rates $\rho_m$ are not known a priori or they are approximately estimated "on the fly," so that $\alpha_m \neq 1/\rho_m$, (6) still represents a useful, measurable quantity for practical purposes. In this case, $\overline{D}(N)$ is the weighted average delay where the weight factor equals $\alpha_m \rho_m$.

## 3. Parameter Separation Design of Throughput-Optimal Scheduling Policies

In this section, we study some general characteristics of throughput-optimal scheduling policies.

**Theorem 1.** *A throughput-optimal policy always allocates the rate vector on the convex hull of the instantaneous rate region.*

*Proof.* If the scheduling policy $\mathcal{P}$ allocates a rate vector $\mathbf{r}^{\mathcal{P}}(\mathbf{h}(n),\mathbf{q}(n))$ in the interior of the convex hull of $\mathcal{C}(\mathbf{h}(n),\overline{P})$, we have

$$\boldsymbol{\mu}^{*T}\mathbf{r}^{\mathcal{P}}(\mathbf{h}(n),\mathbf{q}(n)) < \max_{\widetilde{\mathbf{r}}\in\mathcal{C}(\mathbf{h}(n),\overline{P})}\boldsymbol{\mu}^{*T}\widetilde{\mathbf{r}} \quad (7)$$

for some $\boldsymbol{\mu}^* \in \mathbb{R}_+^M$. Disregarding sets of measure zero and since the policy is independent of time index $n$, the ergodic achievable rate region $\mathcal{C}_{\mathrm{erg}}^{\mathcal{P}}(\overline{P})$ of the policy $\mathcal{P}$ is smaller than $\mathcal{C}_{\mathrm{erg}}(\overline{P})$,

$$\mathcal{C}_{\mathrm{erg}}^{\mathcal{P}}(\overline{P}) := \bigcup_{\|\widetilde{\mathbf{q}}\|\in\mathbb{R}_+^M}\{\widetilde{r}_1,\dots,\widetilde{r}_M : \widetilde{\mathbf{r}} \leq \mathbb{E}\{\mathbf{r}^{\mathcal{P}}(\mathbf{h},\widetilde{\mathbf{q}})\}\} \subset \mathcal{C}_{\mathrm{erg}}(\overline{P}). \quad (8)$$

Thus, the scheduling policy does not achieve the entire ergodic rate region $\mathcal{C}_{\mathrm{erg}}(\overline{P})$ and is not throughput-optimal. $\square$

Theorem 1 is to be understood in the sense that if the rates are not allocated on the convex hull of the instantaneous

rate region, some arrival traffic with $\boldsymbol{\rho} \in \mathcal{C}(\mathbf{h}(n),\overline{P})$ can be constructed so that the queueing system is uniformly transient. Since not all arrival rates with $\boldsymbol{\rho} \in \mathcal{C}(\mathbf{h}(n),\overline{P})$ can be supported, the policy is not throughput-optimal.

Therefore, throughput-optimal scheduling policies can be formulated as an optimization problem,

$$\mathbf{r}^{\mathcal{P}}(\mathbf{h}(n),\mathbf{q}(n)) = \arg\max_{\widetilde{\mathbf{r}}\in\mathcal{C}(\mathbf{h}(n),\overline{P})}\boldsymbol{\mu}^{\mathcal{P}}(\mathbf{h}(n),\mathbf{q}(n))^T\widetilde{\mathbf{r}}, \quad (9)$$

where $\boldsymbol{\mu}^{\mathcal{P}}$ determined by scheduling policy $\mathcal{P}$ is a mapping from the current channel state $\mathbf{h}(n)$ and the queue state $\mathbf{q}(n)$ to the set of weight factors. Generally two mappings $\boldsymbol{\mu}^{\mathcal{P}}$ and $\widetilde{\boldsymbol{\mu}}^{\mathcal{P}}$ lead to the same rate point, where the convex hull has no unique supporting hyperplane. In this case, we define $\boldsymbol{\mu}^{\mathcal{P}}$ to be equivalent to $\widetilde{\boldsymbol{\mu}}^{\mathcal{P}}$ if they lead to the same point. The following theorem presents an important property of the mapping $\boldsymbol{\mu}^{\mathcal{P}}$.

**Theorem 2.** *The mapping $\boldsymbol{\mu}^{\mathcal{P}}$ which characterizes a throughput-optimal scheduling policy is independent of the current fading state $\mathbf{h}(n)$.*

*Proof.* We choose arbitrarily a weight vector $\boldsymbol{\mu}^*$ corresponding to a fixed boundary point $\mathbf{r}^*$ of the ergodic achievable rate region, hence $\boldsymbol{\mu}^*$ is independent of the instantaneous channel state. According to Theorem 1 we have for the channel state $\widehat{\mathbf{h}}$ and the queue state $\widehat{\mathbf{q}}$

$$\mathbf{r}^{\mathcal{P}}(\widehat{\mathbf{h}},\widehat{\mathbf{q}}) = \arg\max_{\mathbf{r}\in\mathcal{C}(\widehat{\mathbf{h}},\overline{P})}(\boldsymbol{\mu}^{\mathcal{P}}(\widehat{\mathbf{h}},\widehat{\mathbf{q}}))^T\mathbf{r}. \quad (10)$$

Thus, for fixed $\widehat{\mathbf{q}}\in\mathbb{R}_+^M$, we have

$$\boldsymbol{\mu}^{*T}\mathbb{E}\{\mathbf{r}^{\mathcal{P}}(\widehat{\mathbf{h}},\widehat{\mathbf{q}})\,|\,\widehat{\mathbf{q}}\} = \mathbb{E}\left\{\boldsymbol{\mu}^{*T}\arg\max_{\mathbf{r}\in\mathcal{C}(\widehat{\mathbf{h}},\overline{P})}(\boldsymbol{\mu}^{\mathcal{P}}(\widehat{\mathbf{h}},\widehat{\mathbf{q}}))^T\mathbf{r}\,|\,\widehat{\mathbf{q}}\right\}$$
$$\leq \mathbb{E}\left\{\max_{\mathbf{r}\in\mathcal{C}(\widehat{\mathbf{h}},\overline{P})}\boldsymbol{\mu}^{*T}\mathbf{r}\right\}$$
$$= \boldsymbol{\mu}^{*T}\mathbf{r}^*. \quad (11)$$

Equality holds if and only if $\boldsymbol{\mu}^{\mathcal{P}}(\widehat{\mathbf{h}},\widehat{\mathbf{q}}) = \boldsymbol{\mu}^*$ and the boundary point is achieved by the scheduler $\mathcal{P}$, otherwise the scheduling policy gives a rate vector in the interior of the ergodic rate region. Therefore, if $\boldsymbol{\mu}^{\mathcal{P}}$ is dependent of the instantaneous channel state, we can choose an arrival process whose mean rate $\boldsymbol{\rho}^*$ fulfills

$$\max_{\widehat{\mathbf{q}}\in\mathbb{R}_+^M}\boldsymbol{\mu}^{*T}\mathbb{E}\{\mathbf{r}^{\mathcal{P}}(\widehat{\mathbf{h}},\widehat{\mathbf{q}})\,|\,\widehat{\mathbf{q}}\} < \boldsymbol{\mu}^{*T}\boldsymbol{\rho}^* < \boldsymbol{\mu}^{*T}\mathbf{r}^*. \quad (12)$$

Define a bounded positive function with $\mathbf{x}\in\mathbb{R}_+^M$

$$V(\mathbf{x}) = 1 - \frac{1}{\boldsymbol{\mu}^{*T}\mathbf{x}+1}, \quad (13)$$

we have

$$\Delta V = \mathbb{E}\{V(\mathbf{q}(n+1)) - V(\mathbf{q}(n))|\mathbf{q}(n)\}$$

$$= \mathbb{E}\left\{ \frac{\boldsymbol{\mu}^{*T}[\mathbf{q}(n)-\mathbf{r}(n)+\mathbf{a}(n)]^+ - \boldsymbol{\mu}^{*T}\mathbf{q}(n)}{(\boldsymbol{\mu}^{*T}[\mathbf{q}(n)-\mathbf{r}(n)+\mathbf{a}(n)]^+ +1)(\boldsymbol{\mu}^{*T}\mathbf{q}(n)+1)} \Big| \mathbf{q}(n) \right\}$$

$$\geqslant \mathbb{E}\left\{ \frac{\boldsymbol{\mu}^{*T}\mathbf{a}(n) - \boldsymbol{\mu}^{*T}\mathbf{r}(n)}{(\boldsymbol{\mu}^{*T}(\mathbf{q}(n)+\mathbf{a}(n))+1)(\boldsymbol{\mu}^{*T}\mathbf{q}(n)+1)} \Big| \mathbf{q}(n) \right\},$$

$$(14)$$

if $\|\mathbf{q}(n)\|$ is sufficiently large. Since the arrival rate is bounded $a_m(n) < c_a < +\infty$, for all $m \in \mathcal{M}$, it holds

$$\Delta V > \mathbb{E}\left\{ \frac{\boldsymbol{\mu}^{*T}\mathbf{a}(n) - \boldsymbol{\mu}^{*T}\mathbf{r}^{\mathscr{P}}(\mathbf{h}(n),\mathbf{q}(n))}{(\boldsymbol{\mu}^{*T}\mathbf{q}(n)+c_a\|\boldsymbol{\mu}^*\|+1)(\boldsymbol{\mu}^{*T}\mathbf{q}(n)+1)} \Big| \mathbf{q}(n) \right\}$$

$$= \frac{\boldsymbol{\mu}^{*T}\boldsymbol{\rho}^* - \boldsymbol{\mu}^{*T}\mathbb{E}\{\mathbf{r}^{\mathscr{P}}(\mathbf{h}(n),\mathbf{q}(n))|\mathbf{q}(n)\}}{(\boldsymbol{\mu}^{*T}\mathbf{q}(n)+c_a\|\boldsymbol{\mu}^*\|+1)(\boldsymbol{\mu}^{*T}\mathbf{q}(n)+1)}$$

$$> 0.$$

$$(15)$$

Since the last inequality holds for any $\mathbf{q}(n)$, according to [23, Theorem 8.4.2] the queueing system is uniformly transient. □

As we introduced in Section 1, cross-layer design usually improves the system performance at the cost of high computational complexity. The optimization problem involves variables and constraints from both PHY and MAC layer. The resources, which can be dynamically adapted, are not only limited to transmit power, but can also be extended to code, frequency, and space according to the applied physical model. At the same time, the scheduler must consider the possible evolution of the queue states in subsequent time slots. However, following the result in Theorem 2, we can define the weight vector $\boldsymbol{\mu}^{\mathscr{P}}(\mathbf{q})$ of a throughput-optimal policy as a function only determined by queue state $\mathbf{q}$. In this way, the classical *cross-layer* optimization problem can be divided into two subproblems: finding the optimal weight vector $\boldsymbol{\mu}^{\mathscr{P}}(\mathbf{q})$ according to the queue states; solving the resource allocation problem (9) with the given weight vector. By the separation of the optimization parameters, the complexity of the optimization problem is largely reduced. Since the second subproblem can be efficiently solved for various physical models, the scheduling design problem reduced to find the optimal weight vector for the optimization problem. Particularly for the considered delay optimization problem, the interface between the two subproblems is the weight factor $\boldsymbol{\mu}^{\mathscr{P}}(\mathbf{q})$. An illustration of the scheme is shown in Figure 2. The average packet delay $\overline{D}(N)$ is dependent on the rate allocation $\mathbf{r}^{\mathscr{P}}$, which is controlled by the weight factor $\boldsymbol{\mu}^{\mathscr{P}}$. Thus in Subproblem 1 we aim to find the optimal weight factor which minimize averaged delay $\overline{D}(N)$. The obtained weight factor $\boldsymbol{\mu}^{\mathscr{P}}$ is then used in Subproblem 2 to calculate the rate allocation $\mathbf{r}$. The details of this scheduling algorithm is introduced in the next section.
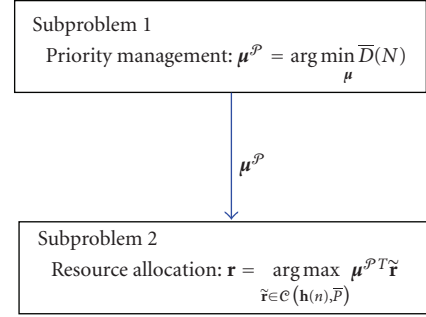


FIGURE 2: Illustration of parameter separation. The suboptimization problems are coupled by the weight factor $\boldsymbol{\mu}^{\mathscr{P}}$.

## 4. Scheduling Design

In this section, we introduce our scheduling policy. First, we solve the delay-optimization problem for a queueing system with a static channel and no new packet arrivals. Then, we adapt the scheduling policy to the queueing system with dynamic channels and random packet arrivals.

*4.1. Scheduling Policy for a Static Channel.* Consider a static channel $\hat{\mathbf{h}}$ and the initial queue states $\mathbf{q}(n=1)$, we assume there is *no* packet arriving after $n = 1$ and we choose a sufficiently large observation time window $N$ with $q_m(N) = 0$, for all $m \in \mathcal{M}$, so that the buffers are completely emptied within the time window. Thus, the scheduling policy can be written as the solution of the optimization problem:

$$\min \sum_{m=1}^{M} D_m(N) \equiv \min \sum_{n=1}^{N} \sum_{m=1}^{M} \alpha_m q_m^n$$

$$\text{s.t.} \quad q_m^{n+1} = q_m^n - r_m^n,$$

$$\mathbf{r}^{\mathbf{n}} \in \mathcal{C}(\hat{\mathbf{h}}, \overline{P}),$$

$$q_m^n - r_m^n \geq 0, \quad \forall m \in \mathcal{M}, n \in [1,\dots,N],$$

$$(16)$$

where $q_m^n$, $r_m^n$ denote the queue length and transmit rate of user $m$ in time slot $n$. For convenience, we also use the superscript to denote the time slot in the following. Extending the problem (16) to each queue state $\mathbf{q}^n$ we have the equivalent optimization problem

$$\min \sum_{n=1}^{N} \left( \sum_{m=1}^{M} \alpha_m q_m^1 - \sum_{m=1}^{M} (N-n)\alpha_m r_m^n \right)$$

$$\text{s.t.} \quad \mathbf{r}^n \in \mathcal{C}(\hat{\mathbf{h}}, \overline{P}), \qquad (17)$$

$$q_m^1 - \sum_{t=1}^{n} r_m^t \geq 0, \quad \forall m \in \mathcal{M}, n \in [1,\dots,N].$$

The problem (17) states a convex optimization problem and we can solve it using standard "ready-to-use" methods. However, this problem involves parameters over $N$ time slots and $M$ users, which is very complicated, especially if $N$ is

(1) Set $\mu_m^{(0)} = \alpha_m$, for all $m \in \mathcal{M}$ and calculate $\mathbf{r}^{(0)} = \arg\max_{\mathbf{r} \in \mathcal{C}(\mathbf{h}, \overline{P})} \boldsymbol{\mu}^{(0)T} \mathbf{r}$.
(2) Initialize the length of non-idle state $\eta_m^{(0)} = \min_{m \in \mathcal{M}}(q_m^1 / r_m^{(0)})$, for all $m \in \mathcal{M}$.
(3) Set the order $\pi$ so that $q_{\pi(1)}^1 / r_{\pi(1)}^{(0)} \geq q_{\pi(2)}^1 / r_{\pi(2)}^{(0)} \geq \cdots \geq q_{\pi(M)}^1 / r_{\pi(M)}^{(0)}$.
(4) Set $t = 0$.
**repeat**
    (5.1) Set $\boldsymbol{\eta}^{(t+1)} = \boldsymbol{\eta}^{(t)}$
  **for** $m = 1$ to $M$ **do**
      (5.2.1) $\boldsymbol{\eta}' = \boldsymbol{\eta}^{(t+1)}$
     **repeat**
        (5.2.2.1) Increase $\eta'_{\pi(m)}$ and calculate $\boldsymbol{\mu}^n$ by setting $\eta^*_{\pi(m)} = \eta'_{\pi(m)}$ in (19).
        (5.2.2.2) Solve the maximization problem (20) and calculate the evolution of the queue state.
        **if** $q_{\pi(m)}^{\lceil \eta'_{\pi(m)} \rceil} \geq 0$ **then**
            $\eta^{(t+1)}_{\pi(m)} = \eta'_{\pi(m)}$
        **end if**
        **until** $q_{\pi(m)}^{\lceil \eta'_{\pi(m)} \rceil} < 0$
    **end for**
    $t = t + 1$
  **until** $\eta_m^{(t)} - \eta_m^{(t-1)} < \epsilon$, for all $m \in \mathcal{M}$
(6) $\boldsymbol{\eta}^* = \boldsymbol{\eta}^{(t)}$

$\epsilon$ is the predefined error tolerance of $\eta_m$.

ALGORITHM 1: Idle state prediction algorithm.

large. Therefore, we introduce in the following an iterative algorithm called *idle state prediction* algorithm to solve the problem.

Formulating the Lagrangian function of (17)

$$
L(\mathbf{r}^n, \boldsymbol{\lambda}^n) = \sum_{n=1}^{N} \sum_{m=1}^{M} \alpha_m q_m^1 - \sum_{n=1}^{N} \sum_{m=1}^{M} (N-n) \alpha_m r_m^n
$$
$$
- \sum_{n=1}^{N} \sum_{m=1}^{M} \lambda_m^n \left( q_m^1 - \sum_{t}^{n} r_m^t \right). \tag{18}
$$

Denote $\eta_m^* = \alpha_m N - \sum_{t=1}^{N} \lambda_m^t$, if $\eta_m^*$ is known, we can get the optimal $\mu_m^n$ with

$$
\mu_m^n = \begin{cases} \alpha_m(\eta_m^* - n + 1) & n \leq \eta_m^*, \\ 0 & n > \eta_m^*, \end{cases} \tag{19}
$$

and the delay-optimization problem is transformed into

$$
\max \sum_{n=1}^{N} \boldsymbol{\mu}^{nT} \mathbf{r}^n \tag{20}
$$
$$
\text{s.t.} \quad \mathbf{r^n} \in \mathcal{C}(\widehat{\mathbf{h}}, \overline{P}),
$$

where $\boldsymbol{\mu}^n$ is the vector of weight factors in the $n$th time slot.

The parameters $\eta_m^*$ in (19) can be interpreted as the expected service time of user $m$ if the optimal solution is applied. In the time slots $n > \eta_m^*$ the buffer of user $m$ is emptied and the corresponding transmitter is in idle state. Based on this property, $\boldsymbol{\eta}^*$ is obtained with an iterative approach given in Algorithm 1.

**Theorem 3.** *$\boldsymbol{\eta}^{(t)}$ obtained in Algorithm 1 converges to the to $\boldsymbol{\eta}^*$ which gives the optimal $\mu_i^n$ for the delay-optimization problem (16).*

*Proof.* In any time slot $n > \eta_m^*$, we have $\mu_m^n = 0$ which means the buffer of $i$th user is empty at the $n$th time slot. In any $n \leq \eta_m^*$, the $m$th buffer must be nonempty. Therefore, if $\boldsymbol{\eta}^{(t)} = \boldsymbol{\eta}^*$, we have $q_m(\lceil \eta_m^* \rceil) = 0$, for all $m \in \mathcal{M}$, and $\boldsymbol{\eta}^{(t+1)} = \boldsymbol{\eta}^* = \boldsymbol{\eta}^{(t)}$. The break condition in Step (6) is fulfilled and the algorithm stops at the optimum.

For two users $i, j \in \mathcal{M}$, if $\eta_i^* \geq \eta_j^*$, the optimal weight factors

$$
\frac{\mu_i^n}{\mu_j^n} = \frac{\alpha_i(\eta_i^* - n + 1)}{\alpha_j(\eta_j^* - n + 1)} > \frac{\alpha_i}{\alpha_j} = \frac{\mu_i^{(0)}}{\mu_j^{(0)}}, \quad n \in [1, \dots, \lceil \eta_j^* \rceil]. \tag{21}
$$

The rate allocation is determined by the given weight factor, thus

$$
\frac{r_i^n}{r_j^n} > \frac{r_i^{(0)}}{r_j^{(0)}}, \quad n \in [1, \dots, \lceil \eta_j^* \rceil] \tag{22}
$$

follows.

From (22), we have

$$
\frac{r_i^{(0)}}{r_j^{(0)}} \leq \frac{\sum_{n=1}^{\lceil \eta_j^* \rceil} r_i^n}{\sum_{n=1}^{\lceil \eta_j^* \rceil} r_j^n} \leq \frac{\sum_{n=1}^{\lceil \eta_j^* \rceil} r_i^n + \sum_{n=\lceil \eta_j^* \rceil + 1}^{\lceil \eta_i^* \rceil} r_i^n}{\sum_{n=1}^{\lceil \eta_i^* \rceil} r_j^n} = \frac{q_i^1}{q_j^1}, \tag{23}
$$

and it holds

$$
\frac{q_i^1}{r_i^{(0)}} \geq \frac{q_j^1}{r_j^{(0)}}. \tag{24}
$$

Therefore, $\pi$ in Step (3) gives also the order of $\boldsymbol{\eta}^*$ so that $\eta^*_{\pi(i)} \geq \eta^*_{\pi(j)}$, if $i < j$.

For the $\pi(M)$th user with $\eta^*_{\pi(M)} = \min_{i \in \mathcal{M}} \eta^*_{\pi(i)}$, we have

$$\frac{\mu^n_{\pi(M)}}{\mu^n_i} = \frac{\alpha_{\pi(M)}(\eta^*_{\pi(M)} - n + 1)}{\alpha_i(\eta^*_i - n + 1)} \leq \frac{\alpha_{\pi(M)}}{\alpha_i} = \frac{\mu^{(0)}_{\pi(M)}}{\mu^{(0)}_i}, \quad (25)$$

$$\forall n \in [1, \ldots, \lceil \eta^*_{\pi(M)} \rceil], \ \forall i \in \mathcal{M}, \ i \neq \pi(M).$$

It follows $r^n_{\pi(M)} \leq r^{(0)}_{\pi(M)}$, for all $n \in [1, \ldots, \lceil \eta^*_{\pi(M)} \rceil]$. Hence, in the initial state $\eta^{(0)}_{\pi(M)} \leq \eta^*_{\pi(M)}$ and further $\eta^{(0)}_m = \eta^{(0)}_{\pi(M)} \leq \eta^*_m$, for all $m \in \mathcal{M}$, $m \neq \pi(M)$. In each iteration step, if the optimum $\boldsymbol{\eta}^*$ is not achieved, $\eta^{(t+1)}_m$ can always be increased so that $\eta^{(t+1)}_m > \eta^{(t)}_m$, for all $m \in \mathcal{M}$. Hence, the convergence of the algorithm is proven. □

*4.2. Scheduling Policy for Dynamic Channels.* It is worth noting that if channel state $\mathbf{h}(n)$ varies over time and the base station has the knowledge of each channel state in advance, the algorithm in the previous subsection can also be used in this case with some modifications. However, such a noncausal scheduler is not realizable. The base station has usually only the current channel state information and the statistical knowledge of the channel. In this case, the optimal weight factors are calculated by the ergodic achievable rate region and under the assumption that there is no new packets arrival. In the next time slot, the weight factors must be recalculated according to the new queue state.

If no new packet arrives after the time slot $n = 0$, the expected delay for a given policy $\mathcal{P}$ is

$$\frac{1}{M} \mathbb{E}\left\{ \sum_{m=1}^{M} D_m(N) \right\}$$
$$= \frac{1}{M} \mathbb{E}\left\{ \sum_{n=1}^{N} \left( \sum_{m=1}^{M} \alpha_m q^1_m - \sum_{m=1}^{M} (N - n) \alpha_m r^{\mathcal{P}n}_m \right) \right\}, \quad (26)$$

where $r^{\mathcal{P}n}_m$ is the rate allocated by the policy $\mathcal{P}$ for the $m$th user at $n$th time slot. From Theorem 2, we know that if $\mathcal{P}$ is a throughput-optimal policy, then

$$\mathbb{E}\{\mathbf{r}^{\mathcal{P}}(\mathbf{h}^n, \mathbf{q}^n)\} = \arg \max_{\mathbf{r} \in \mathcal{C}_{\mathrm{erg}}(\overline{P})} (\boldsymbol{\mu}^{\mathcal{P}}(\mathbf{q}^n))^T \mathbf{r}, \quad (27)$$

where $\boldsymbol{\mu}^{\mathcal{P}}$ is independent of the current channel state. Hence, the optimization problem is equivalent to

$$\min \sum_{n=1}^{N} \left( \sum_{m=1}^{M} \alpha_m q^1_m - \sum_{m=1}^{M} (N - n) \alpha_m \tilde{r}^n_m \right)$$

$$\text{s.t.} \quad \tilde{\mathbf{r}}^{\mathbf{n}} \in \mathcal{C}_{\mathrm{erg}}(\overline{P}), \quad (28)$$

$$q^1_m - \sum_{t=1}^{n} r^t_m \geq 0, \quad \forall m \in \mathcal{M}, n[\in 1, \ldots, N].$$

Then, the optimization problem can be solved using Algorithm 2.

In the system with new packet arrivals, the weight vector $\tilde{\boldsymbol{\mu}}$ should be recalculated according to the new queue state and the rate allocation is determined by $\tilde{\boldsymbol{\mu}}$ and the current channel state $\mathbf{h}(n)$.

As we introduced in Section 2, if we chose the factor $\alpha_m = 1/\rho_m$, the limit $\lim_{N \to +\infty} \overline{D}(N)$ represents the average delay of each packet. Average arrival rates $\rho_m$ can be estimated by previous arrival processes. However, even if the estimation deviates from the actual arrival rate, $\lim_{N \to +\infty} \overline{D}(N)$ can still be considered as a useful delay measurement.

The ergodic achievable rate region is also estimated based on the history. The ergodic region is calculated from a number of sampled fading states, thus the computational complexity might be very high. In [9], a method is introduced to approximate the boundary surface of $\mathcal{C}_{\mathrm{erg}}(\overline{P})$ by utilizing a hypersphere. Only $M + 1$ boundary points on $\mathcal{C}_{\mathrm{erg}}(\overline{P})$ are necessary to characterize the hypersphere so that the complexity is significantly reduced.

In order to prove the throughput-optimality of the policy, we need some technical propositions. The following propositions show the scheduling behavior as the queue length in the system increases. Supposing that the queue length of some users are bounded by some constant $c \geq 0$, while the sum of queue length $\|\mathbf{q}\|$ is increasing, we denote the set of these users as $\mathcal{G}_1 := \{m \mid q_m \leq c, m \in \mathcal{M}\}$ and the remainder as $\mathcal{G}_2 = \mathcal{M}/\mathcal{G}_1$.

**Proposition 1.** *If $q_i$ is bounded, $i \in \mathcal{G}_1$, and $q_j$ is unbounded, $j \in \mathcal{G}_2$ while $\|\mathbf{q}\|$ is increasing, there exists some $B > 0$, so that*

$$\frac{\eta_i}{\eta_j} < \epsilon_1, \quad \forall \|\mathbf{q}\| > B \quad (29)$$

*for arbitrary $\epsilon_1 > 0$.*

*Proof.* We denote the expected service time for user $i$, $j$ as $\eta_i$, $\eta_j$ and $\eta'_i$, $\eta'_j$ where the initial queue length of user $i$ is fixed to $q_i$ and the queue length of user $j$ is increased such that $q'_j > q_j$.

Suppose

$$\frac{\eta'_i}{\eta'_j} \geq \frac{\eta_i}{\eta_j}, \quad (30)$$

then the weight factor

$$\frac{\mu^{n'}_i}{\mu^{n'}_j} = \frac{\alpha_i(\eta'_i - n + 1)}{\alpha_j(\eta'_j - n + 1)} \geq \frac{\alpha_i(\eta_i - n + 1)}{\alpha_j(\eta_j - n + 1)} = \frac{\mu^n_i}{\mu^n_j}, \quad (31)$$
$$\forall n \in [1, \ldots, \lceil \eta_i \rceil].$$

The rate allocation $r^n_i$ and $r^{n'}_i$ are determined by the weight factor $\mu^n_i$, $\mu^{n'}_i$, then it holds

$$r^{n'}_i \geq r^n_i, \quad \forall n \in [1, \ldots, \lceil \eta_i \rceil]. \quad (32)$$

Further, we have

$$q_i = \sum_{n=1}^{\lceil \eta'_i \rceil} r^{n'}_i > \sum_{n=1}^{\lceil \eta_i \rceil} r^{n'}_i \geq \sum_{n=1}^{\lceil \eta_i \rceil} r^n_i = q_i, \quad (33)$$

---

**for** each time slot $n$, **do**

    (1) Calculate $\eta^*$ according to current queue state $\mathbf{q}$ using Algorithm 1, where the static rate region $\mathcal{C}(\hat{\mathbf{h}}, \overline{P})$ is replaced with the ergodic achievable rate region $\mathcal{C}_{\text{erg}}(\overline{P})$.

    (2) Calculate the weight factor $\tilde{\mu}_m = \alpha_m \eta_m^*$, for all $m \in \mathcal{M}$ for the current time slot.

    (3) Calculate the current rate allocation

$$\mathbf{r}^* = \arg \max_{\mathbf{r} \in \mathcal{C}(\mathbf{h}^n, \overline{P})} (\tilde{\boldsymbol{\mu}})^T \mathbf{r},$$

    where $\mathbf{h}^n$ is the current channel state.

**end for**

---

ALGORITHM 2

and reach a contradiction. Therefore, it is shown that if $q_i$ is fixed, $\eta_i/\eta_j$ monotonously decreases with growing $q_j$ as long as $\eta_i/\eta_j > 0$ and the proof follows. $\qquad\square$

**Proposition 2.** *If $q_i$, $q_j$ are unbounded, $i, j \in \mathcal{G}_2$, while $\|\mathbf{q}\|$ is increasing such that $q_i = \gamma_i \|\mathbf{q}\|$, $q_j = \gamma_j \|\mathbf{q}\|$, for some $\gamma_i, \gamma_j > 0$, there exist some $B > 0$, so that*

$$\left| 1 - \frac{\eta_i}{\eta_j} \right| < \epsilon_2, \quad \forall \|\mathbf{q}\| > B \tag{34}$$

*for arbitrary $\epsilon_2 > 0$.*

*Proof.* Consider two queue state $\mathbf{q}$ and $\mathbf{q}'$, $\|\mathbf{q}'\| = \theta\|\mathbf{q}\|$ for some $\theta > 1$. The estimated service time for queue state $\mathbf{q}$ and $\mathbf{q}'$ is denoted as $\boldsymbol{\eta}$, $\boldsymbol{\eta}'$ and the expected rate allocation at the $n$th time slot is $\mathbf{r}^n$ and $\mathbf{r}^{n'}$. Without loss of generality, we set $\eta_i/\eta_j > 1$. Suppose

$$\frac{\eta_i'}{\eta_j'} \ge \frac{\eta_i}{\eta_j}, \tag{35}$$

it holds

$$\frac{(1/\eta_i')\sum_{n=1}^{\lceil \eta_i' \rceil} r_i^{n'}}{(1/\eta_j')\sum_{n=1}^{\lceil \eta_j' \rceil} r_j^{n'}} = \frac{(1/\eta_i')q_i'}{(1/\eta_j')q_j'} = \frac{(1/\eta_i')\gamma_i}{(1/\eta_j')\gamma_j}$$

$$> \frac{(1/\eta_i)\sum_{n=1}^{\lceil \eta_i \rceil} r_i^n}{(1/\eta_j)\sum_{n=1}^{\lceil \eta_j \rceil} r_j^n} = \frac{(1/\eta_i)q_i}{(1/\eta_j)q_j} = \frac{(1/\eta_i)\gamma_i}{(1/\eta_j)\gamma_j}, \tag{36}$$

then $\eta_i/\eta_j > \eta_i'/\eta_j'$ follows which leads to the contradiction to (35). Hence, $\eta_i/\eta_j$ decreases with $\|\mathbf{q}\|$ as long as $\eta_i/\eta_j > 1$ and the proof follows. $\qquad\square$

**Theorem 4.** *The proposed scheduling policy keeps the system stable for any set of arrival rates of which the expected value $\boldsymbol{\rho}$ lies inside the ergodic achievable rate region.*

*Proof.* For the proof of weak stability it is sufficient to show that for any $\boldsymbol{\rho} \in \mathcal{C}_{\text{erg}}(\overline{P})$, the Lyapunov drift $\Delta V$ is negative for some lower bounded function $V : \mathbb{R}_+^M \to \mathbb{R}_+$ [21, 23]. Supposing that there are $i \in \mathcal{G}_2$ whose queue lengths are unbounded, we choose

$$V(\mathbf{q}) = \sum_{i \in \mathcal{G}_2} \alpha_i q_i. \tag{37}$$

Since $q_i$ is unbounded while $r_i < c_r$, for all $i \in \mathcal{G}_2$, the drift

$$\Delta V = \mathbb{E}\left\{ \sum_{i \in \mathcal{G}_2} \alpha_i [q_i^n + a_i^n - r_i^n]^+ - \sum_{i \in \mathcal{G}_2} \alpha_i q_i^n \,\Big|\, \mathbf{q}^n \right\}$$

$$= \mathbb{E}\left\{ \sum_{i \in \mathcal{G}_2} \alpha_i (a_i^n - r_i^n) \,|\, \mathbf{q}^n \right\}. \tag{38}$$

Choose arbitrary $j \in \mathcal{G}_2$ and according to Propositions 1 and 2, we have

$$\Delta V \le \mathbb{E}\left\{ \sum_{i \in \mathcal{G}_2} \alpha_i \left( \frac{\eta_i}{\eta_j} + \epsilon_2 \right) (a_i^n - r_i^n) \,\Big|\, \mathbf{q}^n \right\}$$

$$\le \frac{1}{\eta_j} \sum_{i \in \mathcal{G}_2} \alpha_i (\eta_i \rho_i - \mathbb{E}\{r_i^n | \mathbf{q}^n\}) + \sum_{i \in \mathcal{G}_2} \epsilon_2 \alpha_i (c_a + c_r)$$

$$\le \frac{1}{\eta_j} \sum_{i \in \mathcal{M}} \alpha_i \eta_i (\rho_i - \mathbb{E}\{r_i^n | \mathbf{q}^n\}) + \sum_{i \in \mathcal{G}_2} \epsilon_2 \alpha_i (c_a + c_r)$$

$$+ \sum_{i \in \mathcal{G}_1} \epsilon_1 \alpha_i (c_a + c_r). \tag{39}$$

Define $\beta = \max_{\mathbf{r} \in \mathcal{C}_{\text{erg}}(\overline{P})} \min_{i \in \mathcal{M}} (r_i - \rho_i)$, we have

$$\Delta V \le \frac{1}{\eta_j} \sum_{i \in \mathcal{M}} \alpha_i \eta_i (-\beta) + \sum_{i \in \mathcal{G}_2} \epsilon_2 \alpha_i (c_a + c_r)$$

$$+ \sum_{i \in \mathcal{G}_1} \epsilon_1 \alpha_i (c_a + c_r) \tag{40}$$

$$\le -\beta \sum_{i \in \mathcal{M}} \alpha_i + \epsilon_2 \beta \sum_{i \in \mathcal{G}_2} \alpha_i + \epsilon_2 \sum_{i \in \mathcal{G}_2} \alpha_i (c_a + c_r)$$

$$+ \epsilon_1 \sum_{i \in \mathcal{G}_1} \alpha_i (c_a + c_r). \tag{41}$$

Since the first addend in (41) is constant and the last three addends vanish by increasing $\|\mathbf{q}^n\|$, the drift $\Delta V < 0$ for $\|\mathbf{q}^n\| > B$ if $B$ is sufficiently large. Hence, the Markov chain is weakly stable and the proof follows. $\qquad\square$

## 5. Numerical Evaluations

In order to evaluate the delay performance of the introduced ISP policy, we compare our policy with two other
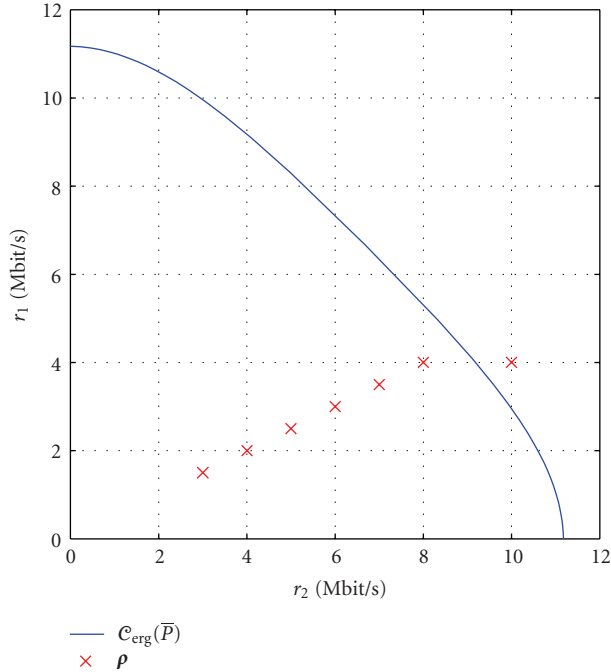
FIGURE 3: Ergodic achievable rate region of an OFDMA system for 2 users. 7 sets of arrival rates (x-marks in the figure) are chosen from the inside/outside of the rate region to test the throughput and delay performance of the system.

throughput-optimal policies: MWMS [3] and QPS [9]. MWMS uses the queue length as the weight factor in the maximization problem:

$$\mathbf{r}^{\mathscr{P}}(\mathbf{h}^n, \mathbf{q}^n) = \arg\max_{\mathbf{r} \in \mathcal{C}(\mathbf{h}^n, \overline{P})} \mathbf{q}^{nT} \mathbf{r}. \qquad (42)$$

For the QPS policy, the weight vector is chosen as the norm at the boundary point of $\mathcal{C}_{\mathrm{erg}}(\overline{P})$, where $\mathbb{E}\{\mathbf{r}^{\mathscr{P}}(\mathbf{h}^n, \mathbf{q}^n)|\mathbf{q}^n\}$ is proportional to $\mathbf{q}^n$,

$$\mathbb{E}\{\mathbf{r}^{\mathscr{P}}(\mathbf{h}^n, \mathbf{q}^n)\} = \mathbf{q}^n \max_{x\mathbf{q}^n \in \mathcal{C}_{\mathrm{erg}}(\overline{P})} x,$$
$$\text{s.t.} \quad \mathbf{r}^{\mathscr{P}}(\mathbf{h}^n, \mathbf{q}^n) \in \mathcal{C}(\mathbf{h}^n, \overline{P}), \qquad (43)$$

where $x$ is a scalar.

The performance of the three schedulers are compared for an OFDMA system as described in [20]. The system has 250 orthogonal subcarriers and an entire bandwidth of 2.5 MHz. The multipath channel is modeled as i.i.d block fading and the length of channel impulse response $L_m = 4$, for all $m \in \mathcal{M}$. The length of a time slot $T$ is 2 milliseconds and in every slot 27 OFDM symbols are transmitted per subcarrier. The modulation is adapted to the different channel states on each subcarrier and can be chosen from QPSK, 16QAM, 64QAM. The source data is coded at rate 2/3, so that the decoding error probability at the receiver is lower than 1e-3. For an average receive SNR of 15 dB the ergodic achievable rate region for two users is shown in Figure 3. Note that the small number of users is
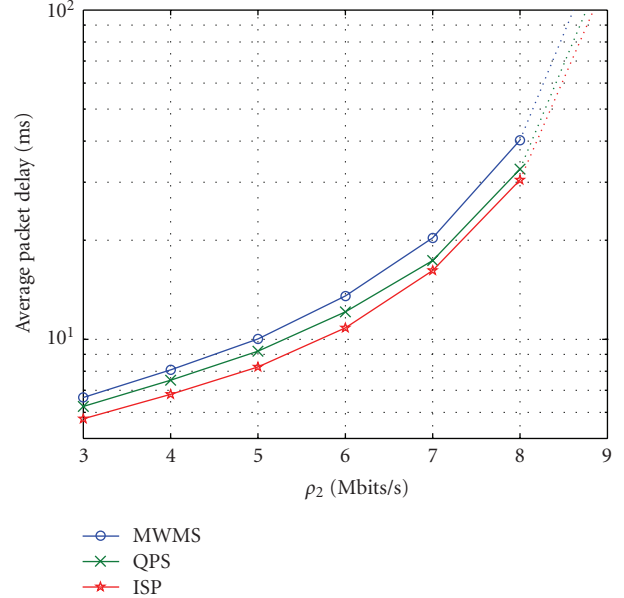


FIGURE 4: Delay performance of the OFDMA system driving by MWMS, QPS, and the introduced ISP scheduling policy. The sets of average arrival rates in the system are indexed by the arrival rate of user 2.

not a limitation but facilitates the description of the ergodic achievable rate region.

We choose arrival rate $\rho_1 = [366, 488, 610, 732, 854, 977, 977]$ packet/s and $\rho_2 = [732, 976, 1220, 1464, 1708, 1954, 2440]$ packet/s for user 1 and 2, respectively, where the size of a packet is 512 bytes. In Figures 3 and 4, the arrival rates are converted to Mbit/s for convenience. In order to verify the stability properties of the system the last set of arrival rates is chosen to lie outside the ergodic achievable rate region.

The average packet delay in the system with the selected sets of arrival rates is shown in Figure 4. It can be seen that for the sets of arrival rates inside the ergodic rate region the system is kept stable in the sense that the average packet delay is finite. For the arrival rates $\rho = [977; 2440]$ packet/s the packet delay tends to go to infinity. All three scheduling policies are throughput-optimal. Compared to the other two scheduling policies, the introduced scheduling policy has the best delay performance and achieves a significant gain. Note that in order to show the rapid growth of the delay time by increased arrival rates, the $y$-axis is logarithmically scaled.

In Figure 5, we compare the delay performance with respect to the number of supported users in the system. The number of users is increased from $M = 2$ to $M = 16$, while the sum of expected arrival rate remains the same. Denote the sum of expected arrival rate by $S_\rho$, we set $\rho_i = 2S_\rho/3M$ for $i \in \mathcal{M}$, if $i$ is odd and $\rho_i = 4S_\rho/3M$ for $i \in \mathcal{M}$, if $i$ is even. The other physical parameters are the same as in Figures 3 and 4. Figure 5 shows the average packet delay in the system resulting from the different scheduling policies. Solid lines are simulated by arrival rate set 1 with $S_\rho = 1098$ packet/s. Dotted lines are simulated by arrival rate set 2 with $S_\rho = 1830$ packet/s. Dashed lines are simulated
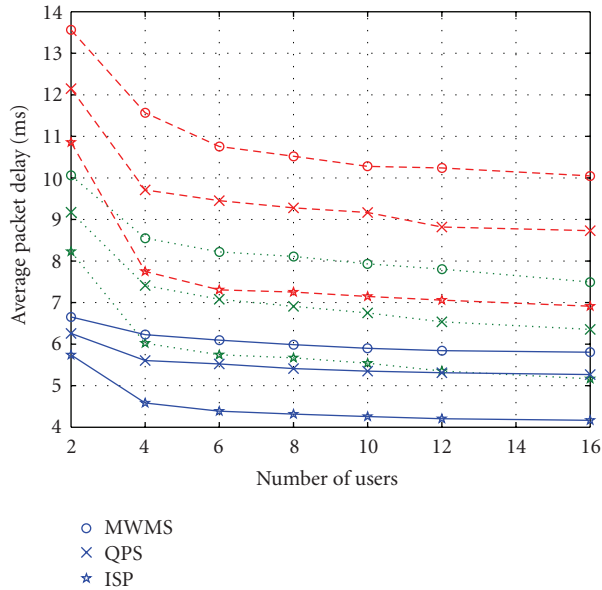
FIGURE 5: Delay performance of MWMS, QPS, and ISP scheduling policy with respect to number of supported users in the system. Solid lines are simulated by arrival rate set 1. Dotted lines are simulated by arrival rate set 2. Dashed lines are simulated by arrival rate set 3.

by arrival rate set 3 with $S_\rho = 2196$ packet/s. Same as in Figure 4, it can be observed that the delay increases with increasing $S_\rho$. Fixing $S_\rho$, the delay decreases with the number of users due to multiuser diversity. At the same time, because of higher flexibility in resource allocation, ISP scheduler provides even more performance gain in delay than the other two schedulers. In case of $M = 6$, ISP scheduler achieves about 30% reduction in averaged delay compared to QPS.

## 6. Conclusion

In this paper, we presented a concept to design throughput-optimal scheduling policies for cellular BC systems. In general it is shown that the scheduling problem can be formulated as a weighted sum rate maximization problem, where the characteristics of the scheduling policy is determined by the choice of weight factors in the maximization problem classifying all throughput-optimal policies. Based on this concept, a throughput-optimal policy is developed to achieve low delay performance. The weight factors achieving the minimum averaged delay are obtained by an iterative procedure, called idle state prediction (ISP) algorithm. The convergence of the algorithm as well as the throughput-optimality of the scheduling policy are proven. Numerical results show that ISP reduces significantly average packet delay compared to other existing scheduling policies. In systems with larger number of users, this advantage becomes even more noticeable due to higher flexibility in resource allocation.

## References

[1] A. Eryilmaz, R. Srikant, and J. R. Perkins, "Stable scheduling policies for fading wireless channels," *IEEE/ACM Transactions on Networking*, vol. 13, no. 2, pp. 411–424, 2005.

[2] G. Wunder and C. Zhou, "Queueing analysis for the OFDMA downlink: throughput regions, delay and exponential backlog bounds," submitted to the *IEEE Transactions on Wireless Communications*.

[3] L. Tassiulas and A. Ephremides, "Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multihop radio networks," *IEEE Transactions on Automatic Control*, vol. 37, no. 12, pp. 1936–1948, 1992.

[4] N. McKeown, A. Mekkittikul, V. Anantharam, and J. Walrand, "Achieving 100% throughput in an input-queued switch," *IEEE Transactions on Communications*, vol. 47, no. 8, pp. 1260–1267, 1999.

[5] M. Andrews, K. Kumaran, K. Ramanan, A. Stolyar, P. Whiting, and R. Vijayakumar, "Providing quality of service over a shared wireless link," *IEEE Communications Magazine*, vol. 39, no. 2, pp. 150–154, 2001.

[6] S. Shakkottai and A. L. Stolyar, "A study of scheduling algorithms for a mixture of real and non-real time data in HDR," Tech. Rep., Bell Laboratories, Murray Hill, NJ, USA, October 2000.

[7] S. Shakkottai and A. L. Stolyar, "Scheduling for multiple flows sharing a time-varying channel: the exponential rule," in *Analytic Methods in Applied Probability: In Memory of Fridrikh Karpelevich*, vol. 207 of *American Mathematical Society Translations, Series 2*, pp. 185–202, American Mathematical Society, Providence, RI, USA, 2002.

[8] A. Eryilmaz, R. Srikant, and J. R. Perkins, "Throughput-optimal scheduling for broadcast channels," in *Modelling and Design of Wireless Networks*, vol. 4531 of *Proceedings of SPIE*, pp. 70–78, Denver, Colo, USA, August 2001.

[9] K. Seong, R. Narasimhan, and J. M. Cioffi, "Queue proportional scheduling via geometric programming in fading broadcast channels," *IEEE Journal on Selected Areas in Communications*, vol. 24, no. 8, pp. 1593–1602, 2006.

[10] M. J. Neely, E. Modiano, and C. E. Rohrs, "Power allocation and routing in multibeam satellites wit time-varying channels," *IEEE/ACM Transactions on Networking*, vol. 11, no. 1, pp. 138–152, 2003.

[11] E. M. Yeh and A. S. Cohen, "Throughput and delay optimal resource allocation in multiaccess fading channels," in *Proceedings of the IEEE International Symposium on Information Theory (ISIT '03)*, p. 245, Yokohama, Japan, June-July 2003.

[12] E. M. Yeh and A. S. Cohen, "Delay optimal rate allocation in multiaccess fading communications," in *Proceedings of the 42nd Allerton Conference on Communication, Control, and Computing*, pp. 140–149, Monticello, Ill, USA, October-September 2004.

[13] W. Lee, K. Seong, and J. M. Cioffi, "Optimal delay region for cross-layer resource allocation," in *Proceedings of the IEEE*

*Global Telecommunications Conference (GLOBECOM '07)*, pp. 3343–3347, Washington, DC, USA, November 2007.

[14] M. Chiang, "Balancing transport and physical layers in wireless multihop networks: jointly optimal congestion control and power control," *IEEE Journal on Selected Areas in Communications*, vol. 23, no. 1, pp. 104–116, 2005.

[15] X. Lin, N. B. Shroff, and R. Srikant, "A tutorial on cross-layer optimization in wireless networks," *IEEE Journal on Selected Areas in Communications*, vol. 24, no. 8, pp. 1452–1463, 2006.

[16] H. Viswanathan, S. Venkatesan, and H. Huang, "Downlink capacity evaluation of cellular networks with known-interference cancellation," *IEEE Journal on Selected Areas in Communications*, vol. 21, no. 5, pp. 802–811, 2003.

[17] R. Böhnke and K.-D. Kammeyer, "Weighted sum rate maximization for the MIMO-downlink using a projected conjugate gradient algorithm," in *Proceedings of the 1st International Workshop on Cross Layer Design (IWCLD '07)*, pp. 82–85, Jinan, China, September 2007.

[18] J. Liu and Y. T. Hou, "Maximum weighted sum rate of multi-antenna broadcast channels," in *Proceedings of the IEEE Global Telecommunications Conference (GLOBECOM '07)*, Washington, DC, USA, November 2007.

[19] L. Li and A. J. Goldsmith, "Capacity and optimal resource allocation for fading broadcast channels—I: ergodic capacity," *IEEE Transactions on Information Theory*, vol. 47, no. 3, pp. 1083–1102, 2001.

[20] G. Wunder, C. Zhou, H.-E. Bakker, and S. Kaminski, "Throughput maximization under rate requirements for the OFDMA downlink channel with limited feedback," *EURASIP Journal on Wireless Communications and Networking*, vol. 2008, Article ID 437921, 14 pages, 2008.

[21] E. Leonardi, M. Mellia, F. Neri, and M. Ajmone Marsan, "Bounds on average delays and queue size averages and variances in input-queued cell-based switches," in *Proceedings of the 20th Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM '01)*, vol. 2, pp. 1095–1103, Anchorage, Alaska, USA, April 2001.

[22] L. Kleinrock, *Queueing Systems*, vol. 1, John Wiley & Sons, New York, NY, USA, 1975.

[23] S. Meyn and R. Tweedie, *Markov Chains and Stochastic Stability*, Springer, London, UK, 1993.