

# REPORT DOCUMENTATION PAGE

Form Approved  
OMB NO. 0704-0188

Public Reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comment regarding this burden estimates or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188,) Washington, DC 20503.

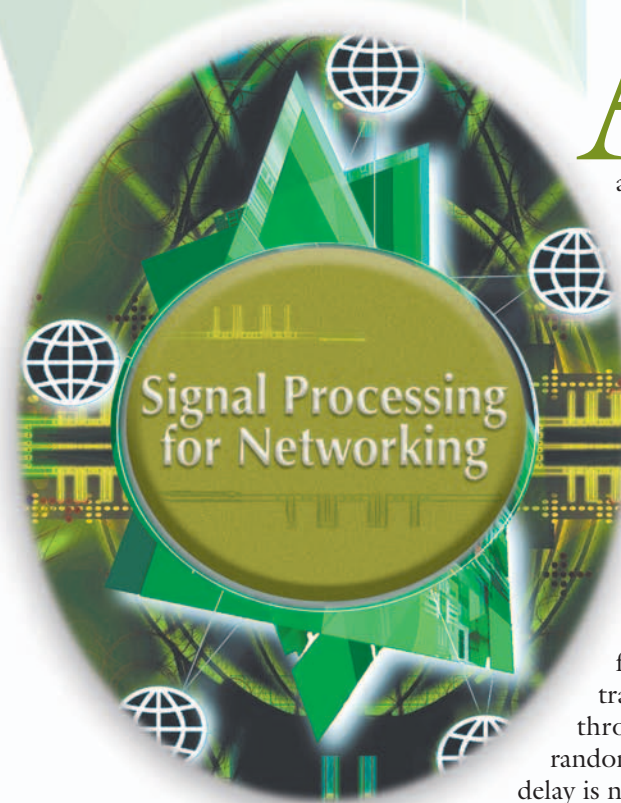
1. AGENCY USE ONLY ( Leave Blank)		2. REPORT DATE June 7, 2005	3. REPORT TYPE AND DATES COVERED	
4. TITLE AND SUBTITLE Cross-layer Wireless Resource Allocation			5. FUNDING NUMBERS DAAD19-03-1-0229	
6. AUTHOR(S) Randall Berry, Dept. of Electrical and Computer Engineering, Northwestern University Edmund Yeh, Dept. of Electrical Engineering, Yale University				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)			8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING / MONITORING AGENCY NAME(S) AND ADDRESS(ES) U.S. Army Research Office P.O. Box 12211 Research Triangle Park, NC 27709-2211			10. SPONSORING / MONITORING AGENCY REPORT NUMBER 45187.5-CI-YIP	
11. SUPPLEMENTARY NOTES The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision, unless so designated by other documentation.				
12 a. DISTRIBUTION / AVAILABILITY STATEMENT Approved for public release; distribution unlimited.			12 b. DISTRIBUTION CODE	
13. ABSTRACT (Maximum 200 words)				
14. SUBJECT TERMS			15. NUMBER OF PAGES	
			16. PRICE CODE	
17. SECURITY CLASSIFICATION OR REPORT UNCLASSIFIED	18. SECURITY CLASSIFICATION ON THIS PAGE UNCLASSIFIED	19. SECURITY CLASSIFICATION OF ABSTRACT UNCLASSIFIED	20. LIMITATION OF ABSTRACT UL	

NSN 7540-01-280-5500

Standard Form 298 (Rev.2-89)  
Prescribed by ANSI Std. Z39-18  
298-102

# Cross-Layer Wireless Resource Allocation

*Randall A. Berry  
and Edmund M. Yeh*



A fundamental problem in networking is the allocation of limited resources among the users of the network. In a traditional layered network architecture, the resource to be allocated at the medium access control (MAC) and network layers is the use of communication links, viewed as “bit pipes” that deliver data at a fixed rate with occasional random errors. This bit pipe is a simple abstraction of the underlying physical and data link layers. This abstraction has, in some ways, caused the research community to split into two distinct groups, which we shall refer to as the networking and communication communities. Research in the networking community has focused on allocating these bit pipes among different streams of randomly arriving traffic using approaches such as packet scheduling and collision resolution. The goal here is to efficiently utilize the bit pipes while providing acceptable quality of service (QoS) in terms of delay and throughput to each user. In contrast, the communication community has focused on building better bit pipes, i.e., improving the transmission rate or spectral efficiency for a given channel through improved detection, modulation, and coding. The random arrivals and departures of traffic are typically ignored and delay is not addressed. Though this separation has many advantages, both practically and conceptually, there is growing awareness that this simple bit-pipe view is inadequate, particularly in the context of modern wireless data networks. Indeed, as highlighted throughout this issue, significant performance gains can be achieved by various *cross-layer* approaches, i.e., approaches that jointly consider physical layer and higher networking layer issues in an integrated framework.

In this article, we consider several basic cross-layer resource allocation problems for wireless fading channels. Here, the resources to be allocated include the transmission

power and rate assigned to each user. In modern wireless systems, a variety of link adaptation techniques such as adaptive modulation and coding or variable-rate spreading are employed that enable a user's data rate to be adapted over time based in part on time-varying channel fading. This results in a physical layer that is no longer well modeled as a fixed-rate bit pipe; instead, a much richer abstraction is required. In this setting, our focus is on characterizing fundamental performance limits, taking into account both network layer QoS and physical layer performance. We note that at the physical layer, fundamental communication limits established by information theory are, in many cases, well understood. However, when higher-layer objectives such as delay are taken into account, much less is known about fundamental performance tradeoffs. The problems surveyed in this article are attempts to address such basic questions. Their solutions serve to establish some benchmarks regarding the achievable performance of cross-layer schemes.

### Cross-Layer Approaches

There are several reasons why cross-layer approaches are particularly well suited for wireless data networks. First, a wireless channel is inherently a shared medium. Efficient resource sharing mechanisms in this setting depend strongly on both the stochastic nature of user activity as well as the selection of physical-layer coding and modulation schemes [1], [2]. For instance, consider a multiaccess problem where a group of distributed users are accessing a common channel. Assuming a simple collision model (i.e., only one user can successfully transmit at any time) leads naturally to the classic ALOHA and carrier sense multiple access (CSMA) algorithms [3], whereas a more code-division multiple access (CDMA)-like model (allowing multiple users to be decoded simultaneously) has very different implications (e.g., [4]). An information-theoretic multiaccess model implies still another set of conclusions [2], [5]–[8]. Second, in wireless networks, where channel quality can vary dramatically in both time and frequency, knowledge of the channel state can be exploited by the system to significantly improve performance. For example, at the physical layer, in a single-user fading channel, the transmission scheme that maximizes the long-term throughput results in transmitting more information in good channel states and less in poor conditions [9]. However, when packet delay is taken into account, it may not be feasible to delay transmission until channel conditions improve. In a multiuser setting, another important characteristic is that channel quality varies across the user population. This results in the phenomenon of *multiuser diversity* [10], whereby as the number of users in a system increases, the probability that some user has a very good channel also increases. Exploiting this diversity results in a total system capacity that is increasing with the number of users. Once again, however, this must be balanced with

network layer issues such as fairness and delay. Finally, the efficient use of energy in mobile devices is of paramount concern in wireless networks. This turns out to be an issue which cuts across almost every protocol layer. In particular, reducing the transmission energy used at the physical layer may result in higher error rates or lower transmission rates, which again affects network layer performance. All of the above coupling effects demonstrate the need to consider network-layer quality of service issues such as throughput and delay jointly with physical-layer issues such as channel fading, coding, and modulation.

We focus primarily on multiaccess (uplink) models, i.e., communication from mobile users to a single base-station or access point. We will also point out several issues that apply to broadcast (downlink) models as well. We consider a situation where randomly arriving data is buffered until it is transmitted and resources are allocated as a function of each flow's buffer occupancy and channel state. We are primarily interested in the case where a centralized controller makes all resource allocation decisions, though some comments about distributed approaches are also included. To characterize fundamental performance limits, we address these problems within an information theoretic framework. Specifically, when allocating resources, such as rate and power, these quantities are constrained by the appropriate capacity region, which depends on the current channel state. Since information theoretic capacity regions characterize asymptotic limits, requiring arbitrary long coding lengths, a careful reader may argue that such results are not applicable in a setting where delays are important. We address this issue in two ways. First, no matter what code lengths are used in practice, information theory provides an upper bound to all achievable rates. For example, for each channel considered here, a corresponding converse coding theorem [11] establishes that reliable communication is impossible outside the capacity region, *for all coding lengths*. Second, the gap between information-theoretic limits and the performance of practical codes with reasonable complexity has narrowed considerably in recent years, due to rapid advances in coding technology. For fading channels, as long as the coherence time is reasonably long, as is the case in typical situations, it is not unreasonable to assume that powerful codes with manageable block lengths can be employed to approximate information-theoretic limits. (For instance, for a user traveling at urban speeds, the coherence time is typically on the order of tens of milliseconds, while the bandwidth is on the order of megahertz, implying that the coherence time corresponds to a coding length of at least several thousand symbols.) Moreover, channel coherence times are typically much smaller than the relevant time-scales at the network layer. Hence, there is no need to consider using shorter codes to further reduce delays. Finally, we note that the framework presented here is quite general and can accommodate

other physical layer models, such as specific coding and modulation schemes.

## Multiaccess Fading Channels

We consider the multiaccess (uplink) wireless communication setting, where multiple transmitters send to a single receiver, in the same time and frequency locality. Consider an  $M$ -user slowly varying, flat-fading Gaussian multiaccess channel with bandwidth  $W$ , given by the model

$$Y(t) = \sum_{i=1}^M \sqrt{H_i(t)} X_i(t) + Z(t). \quad (1)$$

(For a slowly varying channel, the symbol duration  $T_s$  is much smaller than the channel coherence time  $T_{coh}$ , the time interval over which the fading is roughly constant. Flat fading channels are nonfrequency selective, in the sense that the signal bandwidth  $W$  is much smaller than the channel coherence bandwidth  $B_{coh}$ , the band over which fading is roughly constant.) Here,  $H_i(t)$  is the fading process of the  $i$ th user,  $X_i(t)$  is the transmitted signal of the  $i$ th user,  $Z(t)$  is white Gaussian noise with noise density  $N_0/2$ , and  $Y(t)$  is the received signal. Assume that transmitter  $i$  has a long-term average power constraint  $\bar{P}_i$  and a short-term peak power constraint  $\hat{P}_i$ . (The average power constraint may correspond to a battery energy constraint, while the peak power constraint may correspond to a regulatory constraint.) Next, assume that the channel coherence times are sufficiently long as to allow for long code lengths at a fixed joint fading level  $\mathbf{h}$ .

Given that the  $i$ th transmitter experiences a fixed channel fading level  $h_i$  and employs a fixed power level  $p_i$ , the information-theoretic multiaccess capacity region  $\mathcal{C}_{MAC}(\mathbf{h}, \mathbf{p})$  specifies the set of all transmission rates  $\mathbf{r}$  (in bits per second) at which reliable communication is possible *under any coding and modulation scheme*. This capacity region [11] is the set of all nonnegative vectors  $\mathbf{r}$  such that

$$\sum_{i \in S} r_i \leq W \log \left( 1 + \frac{\sum_{i \in S} h_i p_i}{N_0 W} \right) \quad \text{for all } S \subseteq \{1, \dots, M\}. \quad (2)$$

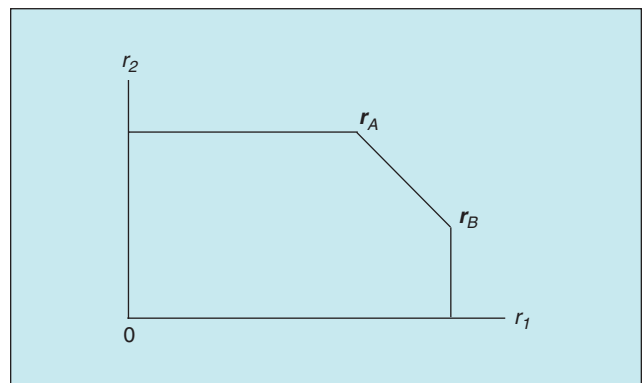
In the two-user case,  $\mathcal{C}_{MAC}(\mathbf{h}, \mathbf{p})$  is a pentagon, as shown in Figure 1. For the  $M$ -user case, it is a bounded convex polyhedron defined by  $2^M - 1$  linear inequalities and  $M$  nonnegativity constraints.

An important observation is that to achieve all rates in the capacity region, joint multi-user coding techniques must be employed. Indeed, CDMA-like strategies, whereby the receiver decodes each user regarding the transmissions of all other users as noise, and simple time-sharing or scheduling strategies, whereby only

one user transmits to the receiver at a time, can typically achieve only a proper subset of the rates in the information-theoretic capacity region (see Figure 1) [1]. To achieve all rates in  $\mathcal{C}_{MAC}(\mathbf{h}, \mathbf{p})$ , a procedure called *successive decoding* can be used. For instance, the corner point  $\mathbf{r}_A$  in Figure 1 is not achievable by pure time-sharing or a CDMA-like strategy, but is achievable by successively decoding user 1 (regarding user 2 as interference in addition to background noise), and then (after subtracting the estimate for user 1 from the received signal), decoding user 2 (facing, with high probability, only background noise). To achieve  $\mathbf{r}_B$ , the receiver implements successive decoding in the opposite order, decoding user 2 first, and then user 1. The successive decoding strategy is generalizable to  $M$  users, and it turns out that  $\mathcal{C}_{MAC}(\mathbf{h}, \mathbf{p})$  has precisely  $M!$  extreme points, one corresponding to each possible permutation of  $\{1, \dots, M\}$  [11].

## Random Arrivals and Resource Allocation

Our focus in this article is on systems where packets for each user arrive to be transmitted at random instants in time. We follow the formulation of [5]–[8]. Specifically, we model the  $i$ th data source as generating packets according to an ergodic counting process  $A_i(t)$ , where  $A_i(t)$  is the number of packet arrivals up to time  $t$ . The packet lengths for source  $i$  are independent identically distributed (i.i.d.) according to distribution function  $F_{Z_i}(\cdot)$  with  $E[Z_i] < \infty$  and  $E[Z_i^2] < \infty$ . Next, assume that each source  $i$ ,  $i = 1, \dots, M$ , has its own (infinite-capacity) buffer into which its packets arrive. Packets for the  $i$ th source are stored in the  $i$ th buffer until they are served by transmitter  $i$ . The transmission power  $P_i(t)$  and rate  $R_i(t)$  used by transmitter  $i$  at time  $t$  are to be dynamically allocated so as to optimize throughput and delay. At the physical layer, we assume that at any time  $t$ , any set of powers and rates from the instantaneous multiaccess information-theoretic capacity region can be allocated to the transmitters, as long as average and peak power constraints are satisfied.

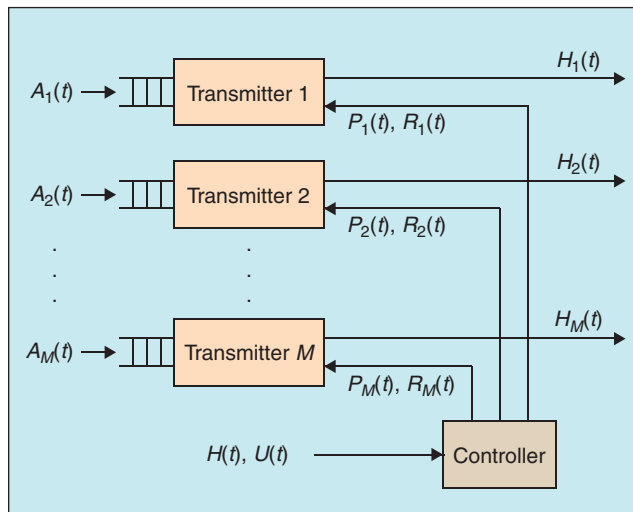


▲ 1. Illustration of  $\mathcal{C}_{MAC}(\mathbf{h}, \mathbf{p})$  for the two-user case. The extreme point  $\mathbf{r}_A$  can be achieved by decoding the users successively in the order (1 2), while  $\mathbf{r}_B$  is achieved by successively decoding in the order (2 1).

We now explicitly pose the dynamic resource allocation problem. Let  $U_i(t)$  be the number of untransmitted bits, or the amount of unfinished work in queue  $i$  at time  $t$ . Consider a stationary controller which at any time  $t \geq 0$  takes as inputs  $\mathbf{H}(t)$  and  $\mathbf{U}(t)$  and outputs a power allocation  $\mathbf{P}(t)$ , and a rate allocation  $\mathbf{R}(t)$ , to transmitters 1 to  $M$ . The controller does this by first choosing a power control policy  $\mathbf{p} = \mathcal{P}(\mathbf{h}, \mathbf{u})$  satisfying  $E[\mathcal{P}_i(\mathbf{H}, \mathbf{U})] \leq \bar{P}_i$  for all  $i$ , where the expectation is taken with respect to the steady-state distribution induced by the controller, and  $\mathcal{P}_i(\mathbf{h}, \mathbf{u}) \leq \hat{P}_i$  for all  $(\mathbf{h}, \mathbf{u})$ , for all  $i$ . Here,  $p_i = \mathcal{P}_i(\mathbf{h}, \mathbf{u})$  is the power allocated to transmitter  $i$  in response to fading state  $\mathbf{h}$  and queue state  $\mathbf{u}$ . Next, the controller chooses a rate allocation policy  $\mathbf{r} = \mathcal{R}(\mathbf{h}, \mathbf{u}) \in \mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathbf{p})$  where  $\mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathbf{p})$  is given by (2). (Note that due to the nature of the constraints, there is no loss of optimality in choosing  $\mathcal{P}$  and  $\mathcal{R}$  in a two-stage manner.) That is, the controller is allowed to adopt stationary power policies  $\mathcal{P}$  that satisfy the average power constraints  $\bar{\mathbf{P}}$  and peak power constraints  $\hat{\mathbf{P}}$ , and given  $\mathcal{P}$ , the controller is allowed to allocate any rate in the multiaccess capacity region induced by the power policy  $\mathcal{P}$ . The setup is illustrated in Figure 2. Our formulation assumes that all transmitters as well as the receiver have access (possibly through side communication channels) to global channel and queue state information  $\mathbf{H}(t)$  and  $\mathbf{U}(t)$ . As described in [12], this setting can often be approximated via feedback in practical wireless systems.

### Stability Region and Throughput Optimal Resource Allocation

The first significant question for the multiaccess queuing system concerns the stability region, i.e., the set all bit arrival rates for which no queue “blows up.” First, some definitions. Let  $\lambda_i = \lim_{t \rightarrow \infty} A_i(t)/t$  denote the packet arrival rate to queue  $i$ , and let  $\rho_i = \lambda_i E[Z_i]$  be the bit arrival rate to queue  $i$ . We define stability as in [13]. Consider the “overflow”



▲ 2. Power and rate allocation for multiaccess fading channels.

function  $f_i(\xi) = \limsup_{t \rightarrow \infty} (1/t) \int_0^t 1_{[U_i(\tau) > \xi]} d\tau$ , where  $1_A$  is the indicator of the event  $A$ . We say that the multiaccess system is *stable* for a particular resource allocation policy if  $f_i(\xi) \rightarrow 0$  as  $\xi \rightarrow \infty$  for all  $i$ . The *stability region* of the multiaccess system is the set of all bit arrival rate vectors  $\boldsymbol{\rho}$  for which there exist some a feasible power control policy and a rate allocation policy under which the system is stable.

It is established in [6] that the stability region is equal to the information-theoretic capacity region under power control, defined in [14]. This region is given by  $\mathcal{C}_{\text{MAC}}(\bar{\mathbf{P}}, \hat{\mathbf{P}}) = \bigcup_{\mathcal{P} \in \mathcal{F}} \mathcal{C}_{\text{MAC}}(\mathcal{P})$  [14]. Here,  $\mathcal{P}$  is a power control policy depending only on the fading state  $\mathbf{h}$  ( $\mathcal{P}(\mathbf{h}, \mathbf{u}) = \mathcal{P}(\mathbf{h})$ ), and  $\mathcal{F}$  is the set of all feasible power control policies depending only on the fading state which satisfy all peak and average power constraints.  $\mathcal{C}_{\text{MAC}}(\mathcal{P})$  is the set of all nonnegative  $\mathbf{r}$  such that  $\sum_{i \in S} r_i \leq E[W \log(1 + (\sum_{i \in S} H_i \mathcal{P}_i(\mathbf{H})/N_0 W))]$  for all  $S \subseteq \{1, \dots, M\}$ . That is,  $\mathcal{C}_{\text{MAC}}(\mathcal{P})$  is the average capacity region (averaged over all fading states) corresponding to a particular power policy  $\mathcal{P} \in \mathcal{F}$ . Suppose joint arrival process  $\{\mathbf{A}(t)\}$  and joint fading process  $\{\mathbf{H}(t)\}$  are modulated by a finite-state ergodic Markov chain. Then, it is shown in [6] that the multiaccess queuing system can be stabilized by some power control and rate allocation policy if  $\boldsymbol{\rho} \in \text{int}(\mathcal{C}_{\text{MAC}}(\bar{\mathbf{P}}, \hat{\mathbf{P}}))$ . Conversely, the multiaccess fading channel cannot be stabilized if  $\boldsymbol{\rho} \notin \mathcal{C}_{\text{MAC}}(\bar{\mathbf{P}}, \hat{\mathbf{P}})$ , as long as the average and peak power constraints are satisfied.

The stability result states that if  $\boldsymbol{\rho} \in \text{int}(\mathcal{C}_{\text{MAC}}(\bar{\mathbf{P}}, \hat{\mathbf{P}}))$ , then the queues can be stabilized. In general, however, this may require knowing the value of  $\boldsymbol{\rho}$ . In reality, the arrival rates  $\boldsymbol{\rho}$  can be learned only over time. One would prefer to find adaptive resource allocation policies which can stabilize the system *without* knowing  $\boldsymbol{\rho}$ , as long as  $\boldsymbol{\rho} \in \text{int}(\mathcal{C}_{\text{MAC}}(\bar{\mathbf{P}}, \hat{\mathbf{P}}))$ , i.e., the system is stabilizable. Such a resource allocation policy is referred to as throughput optimal. Throughput optimal scheduling for fading channels has been examined in [13] and [15]–[19]. These papers, while offering many valuable insights, do not consider information-theoretic optimal coding at the physical layer and do not account for the effect of power control subject to long-term constraints. These important considerations are taken into account in [7], where it is shown that an adaptive version of the power and rate allocation algorithm derived by Tse and Hanly [14] is throughput optimal for the multiaccess queuing system.

In [14], Tse and Hanly consider the problem of maximizing a weighted combination of long-term transmission rates in a multiaccess channel where all transmitters have *infinite backlogs of bits*, and both the transmitters and receivers have access to the channel state. This problem can be stated as

$$\max \sum_{i=1}^M \mu_i r_i \quad \text{subject to } \mathbf{r} \in \mathcal{C}_{\text{MAC}}(\bar{\mathbf{P}}, \hat{\mathbf{P}}) \quad (3)$$

where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_M)$  is a vector of nonnegative weights. Using a Lagrangian formulation of (3) and the underlying polymatroidal structure of  $\mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathbf{p})$ , they show that (3) is equivalent to solving a family of optimization problems over parallel Gaussian multiple access channels, one for each fading state  $\mathbf{h}$ . Their analysis yields a feasible power control policy (satisfying peak and average power constraints) and a rate allocation policy (satisfying capacity constraints) which solve (3). Notice that for a given direction  $\boldsymbol{\mu}$ , the Tse-Hanly power control policy  $\mathcal{P}_{\text{TH}}(\mathbf{h}, \boldsymbol{\mu})$  and rate allocation policy  $\mathcal{R}_{\text{TH}}(\mathbf{h}, \boldsymbol{\mu})$  are *functions of  $\mathbf{h}$  only*.

In [7], it is proved that a throughput optimal resource allocation policy for the multiaccess system with *random packet arrivals* is given by the Tse-Hanly (TH) solution, with the direction  $\boldsymbol{\mu}$  chosen to correspond to the queue state  $\mathbf{u}$ . Specifically, the throughput optimal policy is given by

$$\begin{aligned} \mathcal{P}_{\text{MAC}}^*(\mathbf{h}, \mathbf{u}) &= \mathcal{P}_{\text{TH}}(\mathbf{h}, \boldsymbol{\alpha} * \mathbf{u}), \\ \mathcal{R}_{\text{MAC}}^*(\mathbf{h}, \mathbf{u}) &= \mathcal{R}_{\text{TH}}(\mathbf{h}, \boldsymbol{\alpha} * \mathbf{u}) \end{aligned} \quad (4)$$

where  $\mathbf{u}$  is the queue state,  $\boldsymbol{\alpha}$  is any vector of positive numbers, and  $\boldsymbol{\alpha} * \mathbf{u}$  is the vector whose  $i$ th component is  $\alpha_i u_i$ . The vector  $\boldsymbol{\alpha}$  can be seen as a set of weights representing the relative priorities of the various users. The proof of the throughput optimality results in [7] makes use of the Foster-Lyapunov criterion for the stability of Markov chains [13].

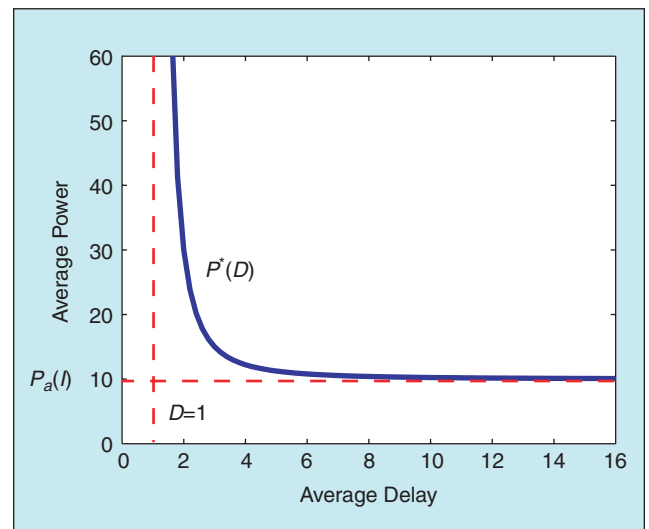
To interpret the throughput optimal resource allocation policies  $\mathcal{P}_{\text{MAC}}^*$  and  $\mathcal{R}_{\text{MAC}}^*$ , let  $\mathbf{v} = \boldsymbol{\alpha} * \mathbf{u}$ , where  $v_i = \alpha_i u_i$ , be the vector of weighted queue sizes. In the case of one user ( $M = 1$ ), it can be shown [14] that  $(\mathcal{P}_{\text{MAC}}^*, \mathcal{R}_{\text{MAC}}^*)$  reduces to the well-known water-filling scheme [9], whereby more power is allocated to favorable channel states, and less or no power is allocated to unfavorable channel states. In the case of multiple users ( $M > 1$ ) where all weighted queue sizes are the same:  $v_1 = v_2 = \dots = v_M$ , and the fading conditions are symmetric,  $(\mathcal{P}_{\text{MAC}}^*, \mathcal{R}_{\text{MAC}}^*)$  reduces to the “multi-user waterfilling” scheme of Knopp and Humblet [10], whereby when all channel states are sufficiently unfavorable, no one transmits. Otherwise, only the user with the best channel condition transmits. In the general case of many users and unequal weighted queue lengths, more than one user typically transmit. Little in general can be said about the optimal power policy  $\mathcal{P}_{\text{MAC}}^*$ . The optimal rate allocation policy  $\mathcal{R}_{\text{MAC}}^*$ , however, satisfies a general principle we refer to as longest weighted queue highest possible rate (LWQHPR). This principle holds for any given feasible power control policy  $\mathcal{P}$ , and is described as follows. Given  $\mathcal{P}$ ,  $\mathcal{R}_{\text{MAC}}^*(\mathbf{h}, \mathcal{P}(\mathbf{h}, \mathbf{u}), \mathbf{u})$  is given by maximizing  $\sum_i v_i r_i$  over  $\mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathcal{P}(\mathbf{h}, \mathbf{u}))$ . Due to the polymatroidal nature of  $\mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathcal{P}(\mathbf{h}, \mathbf{u}))$  [14], the solution is explicitly given as follows. Let  $v_{[1]} \geq v_{[2]} \geq \dots \geq v_{[M]}$  denote the components of  $\mathbf{v}$

in *decreasing* order. Then,  $\mathbf{r}^* = \mathcal{R}_{\text{MAC}}^*(\mathbf{h}, \mathcal{P}(\mathbf{h}, \mathbf{u}), \mathbf{u})$  is given by

$$r_{[i]}^* = W \log \left( 1 + \frac{h_{[i]} \mathcal{P}_{[i]}(\mathbf{h}, \mathbf{u})}{\sum_{j < i} h_{[j]}(t) \mathcal{P}_{[j]}(\mathbf{h}, \mathbf{u}) + N_0 W} \right). \quad (5)$$

It can be verified that  $\mathbf{r}^*$  is the extreme point of  $\mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathcal{P}(\mathbf{h}, \mathbf{u}))$  corresponding to successively decoding the users in the order  $[M], [M - 1], \dots, [1]$ . That is, the smallest component of  $\mathbf{v}$  being decoded first, and the largest component of  $\mathbf{v}$  being decoded last. Alternatively,  $\mathbf{r}^*$  is given by a greedy rate allocation procedure where longer weighted queues are given preference over shorter weighted queues. (Note that the order of decoding is the opposite of the order of preference.) Hence, the name LWQHPR [7]. To illustrate the LWQHPR policy, we refer to the two-user example of Figure 1, where LWQHPR assigns  $\mathbf{r}_A$  (corresponding to decoding order (1 2) if  $v_1 < v_2$  and assigns  $\mathbf{r}_B$  [corresponding to decoding order (2 1)] if  $v_1 \geq v_2$ .

We now compare the performance of the throughput optimal policies  $(\mathcal{P}_{\text{MAC}}^*, \mathcal{R}_{\text{MAC}}^*)$  to those of widely used alternative resource allocation policies. Consider an example in which there are two users observing i.i.d. fading processes. For each user,  $\Pr(H = h_0) = \Pr(H = 1/h_0)$  for some fixed  $h_0 > 0$ . The fading state remains constant for a period of  $T$  seconds and then changes to a new independent fading state. The arrival processes are independent Poisson with  $\lambda_1 = \lambda_2 = \lambda$  and packets lengths are i.i.d. exponential with parameter 1. We focus on the parameters  $h_0 = 10$ ,  $\bar{P} = \bar{P}_1 = \bar{P}_2 = 1$ ,<sup>6</sup>  $T = 0.4$ ,  $N_0 W = 1$ , and equal weights ( $\alpha_1 = \alpha_2$ ). (For simplicity, we assume the peak power constraints are large enough to be ineffective.) Figure 3 shows the simulated performance of the



▲ 3. Total average queue size versus arrival rate for the multiaccess fading channel under five control strategies.

throughput optimal strategy  $(\mathcal{P}_{\text{MAC}}^*, \mathcal{R}_{\text{MAC}}^*)$  relative to those of four other strategies. The sum of the average queue sizes is plotted versus the arrival rate  $\lambda$  for the five strategies described below. The throughput optimal strategy is given by  $(\mathcal{P}_{\text{MAC}}^*, \mathcal{R}_{\text{MAC}}^*)$ . The strategy of Knopp-Humblet [10] maximizes the sum rate assuming an infinite backlog, which corresponds to the throughput optimal strategy with  $\mu_1 = \mu_2$ . The scheduling algorithm allocates power  $2\bar{P}$  to the user with the better fade and zero power to the other user. The constant power longest queue highest possible rate (LQHPR) strategy uses constant power  $(\mathcal{P}_i(\mathbf{h}, \mathbf{u}) = \bar{P}$  for all  $i, \mathbf{h}$  and  $\mathbf{u}$ ) and allocates rates according to (5). The constant power best channel the highest possible rate (BCHPR) strategy also uses constant power, and (ignoring the queue size) gives the BCHPR. The experimental results demonstrate the superior performance of the throughput optimal resource allocation policy, in that its total average queue size is considerably smaller than those of the competitors.

### Delay Optimal Resource Allocation

We have thus far concentrated on stability and throughput optimality. Stability in a queuing system implies that the queue sizes do not “blow up,” but it does not indicate how large the queue sizes can be. To minimize the average packet delay/latency and other related QoS measures, it is necessary to keep the queue sizes *as short as possible*. The general problem of finding delay optimal joint power control and rate allocation policies to minimize average delay for multiaccess channels is still open. In [16], the problem of finding the delay optimal rate allocation policy for a given power control policy is addressed. The main result is that in a symmetric multiaccess queuing system, the symmetric version of the LWQHPR rate allocation policy given by (5) (with  $\alpha_i = 1$  for all  $i$ ) minimizes the average packet delay in a very strong sense.

Consider the case where all arrival processes are non-homogeneous Poisson with rate function  $\lambda(t)$ , and all arriving packets are i.i.d. exponential with common parameter  $\mu > 0$ . Due to the memoryless nature of the system, the vector  $\mathbf{Q}(t) = (Q_1(t), \dots, Q_M(t))$ , where  $Q_i(t)$  is the number of packets in queue  $i$  at time  $t$ , constitutes a state. Thus, we focus on resource allocation policies of the form  $\mathcal{P}(\mathbf{h}, \mathbf{q})$  and  $\mathcal{R}(\mathbf{h}, \mathbf{q})$ , where  $\mathbf{q} = (q_1, \dots, q_M)$  is the vector of queue lengths. We assume that the fading process  $\mathbf{H}(t)$  is symmetric (or exchangeable) in the following sense: for all  $t$ , and all  $\mathbf{a}$  in the fading state space  $\mathcal{H}$ ,  $\Pr(H_1(t) = a_1, \dots, H_M(t) = a_M) = \Pr(H_1(t) = a_{\pi(1)}, \dots, H_M(t) = a_{\pi(M)})$  for any permutation  $\pi$  on the set  $\{1, \dots, M\}$ . Beyond this symmetry, we do not make any other assumptions on the fading process. We focus on power policies  $\mathcal{P}$  which are functions of the fading state only, i.e.,  $\mathcal{P}(\mathbf{h}, \mathbf{q}) = \mathcal{P}(\mathbf{h})$ . We say a power policy is *symmetric* if for all  $\mathbf{a} \in \mathcal{H}$ ,  $\mathcal{P}_i(a_1, \dots, a_M) = \mathcal{P}_{\pi^{-1}(i)}(a_{\pi(1)}, \dots, a_{\pi(M)})$  for any permutation  $\pi$ . That is, under a symmetric

power control policy, the power allocated to a given user is determined by the fading level experienced by that user (and not on the identity of that user) relative to the fading levels experienced by all other queues. For instance, suppose  $M = 2$  and  $a_1 > a_2$ , then if  $\mathcal{P}$  is symmetric,  $\mathcal{P}_1(a_1, a_2) = \mathcal{P}_2(a_2, a_1)$ . An example of a symmetric power control policy is the “multiuser water-filling” policy given by Knopp and Humblet [10].

Consider the version of the LWQHPR policy where  $\alpha_i = 1$  for all  $i$ . We refer to this as the LQHPR policy. The main result on delay optimality from [6] is the following. For an  $M$ -user symmetric multiaccess queuing system, let  $\mathcal{P} : \mathcal{H} \mapsto \mathbb{R}_+^M$  be a given symmetric power control policy. Let  $\mathbf{q}_0$  be the vector of queue sizes at time 0. Let  $\mathbf{Q}(t)$  be the queue evolution under the LQHPR rate allocation policy, and  $\mathbf{Q}'(t)$  be the queue evolution under any feasible rate allocation policy. Then,  $E[\varphi(\mathbf{Q}(t))] \leq E[\varphi(\mathbf{Q}'(t))]$  for all  $t \geq 0$ , for all *increasing and Schur-convex* functions  $\varphi : \mathbb{R}^M \mapsto \mathbb{R}$ . (Schur convex functions are functions which preserve an ordering called majorization [20].) As a main example, the result holds for all symmetric increasing, convex functions on  $\mathbb{R}^M$ . (A function  $\varphi$  is symmetric on  $\mathcal{A} \subset \mathbb{R}^M$  if  $\varphi(\mathbf{x}) = \varphi(\mathbf{x}\mathbf{P})$  for any  $\mathbf{x} \in \mathcal{A}$  and any  $M$  by  $M$  permutation matrix  $\mathbf{P}$ .) More specific examples include  $\varphi(\mathbf{x}) = \max_{i_1 < i_2 < \dots < i_k} (|x_{i_1}| + \dots + |x_{i_k}|)$  for  $1 \leq k \leq M$ ,  $\varphi(\mathbf{x}) = \sum_{i=1}^M |x_i|^r$  for  $r \geq 1$  or  $r \leq 0$ , and  $\varphi(\mathbf{x}) = (\sum_{i=1}^M |x_i|^r)^{1/r}$  for  $r \geq 1$ . Thus, the LQHPR policy minimizes the expected value of a large class of functions of  $\mathbf{Q}(t)$ , of which  $E[\sum_i Q_i(t)]$  is one example. By Little’s law, this implies that LQHPR minimizes the average system delay of packets. The main technique used in proving the result of [6] earlier related results in [5] is *stochastic coupling*, a method relying directly on probabilistic intuition which is capable of generating powerful and elegant results when given certain symmetry conditions. Finally, the delay optimality of the LQHPR policy has been partially extended to the case where arriving packets are i.i.d. according to a general distribution function  $F_Z(\cdot)$  with finite first and second moments [8]. This investigation requires a new analytical technique combining finite-horizon dynamic programming and renewal process theory.

### Broadcast Fading Channels

The analytical framework established earlier for multiaccess networks can be extended to broadcast (down-link) wireless networks, where a single transmitter sends separate *independent* information to  $M$  receivers, where the  $i$ th receiver has fading process  $H_i(t)$  and receiver noise power density  $N_{0i}/2$ . For a fixed transmit power  $p$  and fading states  $h_1, \dots, h_M$ , the broadcast capacity region  $\mathcal{C}_{\text{BC}}(\mathbf{h}, p)$  [11] is very different from the multiaccess region  $\mathcal{C}_{\text{MAC}}(\mathbf{h}, p)$ . As in the multiple-access case, however, simple time-sharing (whereby the transmitter sends to one receiver at a time) cannot achieve all points in  $\mathcal{C}_{\text{BC}}(\mathbf{h}, p)$ . To accomplish the latter, a process called superposition coding com-

bined with the successive decoding technique discussed earlier can be used [11].

In [7], a broadcast communication system with random packet arrivals is considered. It is shown that the queuing system considered in Figure 2 can be directly carried over to the broadcast case, with the following caveats. First, unlike the multiple-access case, there is only one actual transmitter in a broadcast network with long-term average power constraint  $\bar{P}$  and short-term peak power constraint  $\hat{P}$ . There are, however, still multiple arrival processes and queues at the transmitter, corresponding to the information streams of the respective receivers. Since the transmitter uses superposition coding and allocates separate powers and rates to sub-codes for the various receivers, we can associate a “virtual transmitter” with each of the  $M$  receivers (sub-codes). Second, in the broadcast case, the power control policy  $\mathcal{P}$  determines the total power  $p$  used by the actual transmitter as a function of  $\mathbf{h}$  and  $\mathbf{u}$  to satisfy  $E[\mathcal{P}(\mathbf{H}, \mathbf{U})] \leq \bar{P}$ . The rate allocation policy then decides what fraction of power  $\gamma_i p$  (with  $\sum_{i=1}^M \gamma_i = 1$ ) and rate  $r_i$  to assign to each virtual transmitter or subcode.

For the broadcast network, we consider the same throughput and delay questions as in the multiaccess case. The results for stability and throughput optimality are the exact analogues of those for multiaccess channels. In [7], it is shown that stability region for the broadcast network is the same as the information-theoretic capacity region  $\mathcal{C}_{\text{BC}}(\bar{P}, \hat{P})$  under power control defined in [21] and [22]. Moreover, parallel to the multiaccess case, an adaptive version of the power and rate allocation policies designed to maximize transmission rates in a broadcast channel with infinite backlogs of bits, is throughput optimal for a system with random packet arrivals. In [21] and [22], Tse, Li, and Goldsmith use a Lagrangian formulation similar to the multiaccess case to obtain the optimal power policy  $\mathcal{P}_{\text{TLG}}(\mathbf{h}, \boldsymbol{\mu})$  and rate policy  $\mathcal{R}_{\text{TLG}}(\mathbf{h}, \boldsymbol{\mu})$  which maximizes  $\sum_i \mu_i r_i$  subject to  $\mathbf{r} \in \mathcal{C}_{\text{BC}}(\bar{P}, \hat{P})$ , where  $\boldsymbol{\mu}$  is a vector of nonnegative weights. (Even though the structure of  $\mathcal{C}_{\text{BC}}(\mathbf{h}, p)$  and  $\mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathbf{p})$  are very different, greedy algorithms can be used to solve the family of optimization problems, one for each fading state  $\mathbf{h}$ , in both cases.) In [7], it is shown that a throughput optimal policy for the broadcast channel with random packet arrivals is  $\mathcal{P}_{\text{BC}}^*(\mathbf{h}, \mathbf{u}) = \mathcal{P}_{\text{TLG}}(\mathbf{h}, \boldsymbol{\alpha} * \mathbf{u})$ ,  $\mathcal{R}_{\text{BC}}^*(\mathbf{h}, \mathbf{u}) = \mathcal{R}_{\text{TLG}}(\mathbf{h}, \boldsymbol{\alpha} * \mathbf{u})$ , where  $\boldsymbol{\alpha}$  is any vector of positive numbers,  $\boldsymbol{\alpha} * \mathbf{u}$  is the vector whose  $i$ th component is  $\alpha_i u_i$ , and  $u_i$  is the number of untransmitted bits in the queue for the  $i$ th receiver. An interesting consequence of this policy is the following: if there exists a user  $i$  such that  $\alpha_i u_i \geq \alpha_j u_j$  for all  $j \neq i$ , and such that  $(\alpha_i u_i h_i / N_{0i}) \geq (\alpha_j u_j h_j / N_{0j})$  for all  $j \neq i$ , then the throughput optimal policy transmits only to user  $i$ . In particular, if  $\alpha_j = 1$  and  $N_{0j} = N_0$  for all  $j$ , then the throughput optimal policy transmits to the user with best fading when all queues are equal, and transmits to the user with longest queue when all fading levels are equal [7]. Numerical experiments

indicate that the performance of the policy  $(\mathcal{P}_{\text{BC}}^*(\mathbf{h}, \mathbf{u}), \mathcal{R}_{\text{BC}}^*(\mathbf{h}, \mathbf{u}))$  is superior to those of competing resource allocation policies, in that the resulting total average queue size is substantially lower [7].

Finally, we come to the problem of delay optimal resource allocation for broadcast channels. Here, very little progress has been made. Even the problem of finding the delay optimal rate allocation policy for a given power control policy has not been successfully tackled. The delay problem for broadcast channels appears to be more difficult than that for multiaccess channels, mainly because the region  $\mathcal{C}_{\text{BC}}(\mathbf{h}, p)$  lacks the many desirable structural properties of  $\mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathbf{p})$ . Much more work is needed in this area.

## Energy/Delay Tradeoffs

Energy efficiency is a key concern for mobile wireless devices that must rely on limited battery resources. As transmission power is one of the main energy consumers in wireless devices, there has been much interest in approaches for efficiently utilizing this resource. In this context, a fundamental metric is the average energy per bit used in communication. At the physical layer, the minimum energy per bit needed to reliably communicate at a given rate  $R$  can be related to the channel’s capacity. Specifically, let  $C(P)$  denote the capacity of a single user channel as a function of the transmission power. Since for reliable communication  $R < C(P)$ , it follows that the energy per bit,  $E_b$  must satisfy  $E_b > (C^{-1}(R)/R)$ , where  $C^{-1}$  denotes the inverse of  $C(\cdot)$ . For any channel  $(C^{-1}(R)/R)$  is a decreasing convex function of  $R$ ; hence, energy can be conserved by transmitting at a lower rate. For example, if  $C(P) = \ln(1 + P)$  as in a Gaussian noise channel, then  $(C^{-1}(R)/R) = (1/R)(e^R - 1)$ . Asymptotically, as  $R \rightarrow 0$ , the minimum energy per bit is the reciprocal of a channel’s *capacity per unit cost* [23], with “cost” given by energy per channel use.

Transmitting at a lower rate conserves energy, but increases delay at the network layer. This illustrates a fundamental tradeoff between energy efficiency and delay. In fading channels, reducing packet delay may further increase the required energy by forcing users to transmit when channel conditions are not favorable. The energy-delay tradeoff necessitates resource allocation algorithms which optimally balance these two important concerns. A number of approaches have been applied to this problem including [24]–[27]. Interestingly, these approaches are useful even in a single user setting [24], [27], and even in a setting with random arrivals but no channel fading [25].

## Optimal Power/Delay Resource Allocation—Single User Case

We first discuss the case where a single user is transmitting data over the fading channel in (1) (with  $M = 1$ ). Once again, we assume that data randomly arrives and is buffered until it is transmitted. For simplicity, we



consider a discrete-time *block fading* model where during the  $n$ th time-slot the channel gain is constant with value  $H[n]$ , i.e.,  $H(t) = H[n]$  for  $(n-1)\Delta \leq t < n\Delta$  where  $\Delta$  is the length of a time-slot. Denote the unfinished work in the buffer at the start of the  $n$ th time-slot by  $U[n]$  and let  $I[n]$  be the number of bits (amount of work) arriving during the  $n$ th time slot. The unfinished work evolves according to

$$U[n+1] = \max(U[n] + I[n] - R[n]\Delta, 0), \quad (6)$$

where  $R[n]$  is the transmission rate during the  $n$ th time slot. At each time  $n$ , the rate is again specified by a stationary rate allocation policy  $r = \mathcal{R}(h, u)$  that depends on the current queue state and fading state. In this section, we view the system as first specifying a rate allocation policy and then incurring a “power cost”  $P(R[n], H[n])$  that depends on this policy and the current channel state. Specifically, let  $P(r, h)$  represent the minimum power required such that the rate  $r$  is less than the corresponding channel capacity. For the channel in (1), this is given by

$$P(r, h) = \frac{N_0 W}{h} (2^{r/W} - 1), \quad (7)$$

which is an increasing convex function of the transmission rate.

For a given policy  $\mathcal{R}$ , let  $\bar{P}(\mathcal{R}) = \lim_{n \rightarrow \infty} E[P(\mathcal{R}(H[n], U[n]), H[n])]$  denote the steady-state average transmission power consumed. Assuming the system is stable, the average energy per bit is given by  $\bar{P}(\mathcal{R})/\bar{I}$ , where  $\bar{I}$  indicates the average bit arrival rate per time-slot. (Using our previous notation,  $\bar{I} = \lambda \Delta E[Z]$ , where  $E[Z]$  is the expected length of a packet.) Also, for a given policy  $\mathcal{R}$ , let  $\bar{D}(\mathcal{R}) = \lim_{n \rightarrow \infty} E[U[n]/\bar{I}]$  indicate the steady-state average queuing delay. We define the *optimal power/delay tradeoff curve*,  $P^*(D)$ , as

$$P^*(D) = \inf \{ \bar{P}(\mathcal{R}) | \mathcal{R}(\cdot) \text{ such that } \bar{D}(\mathcal{R}) \leq D \}, \quad (8)$$

i.e.,  $P^*(D)$  characterizes the minimum average power required for any policy with an average delay no greater than  $D$ . Conversely, we also define the delay/power tradeoff:

$$D^*(P) = \inf \{ \bar{D}(\mathcal{R}) | \mathcal{R}(\cdot) \text{ such that } \bar{P}(\mathcal{R}) \leq P \}. \quad (9)$$

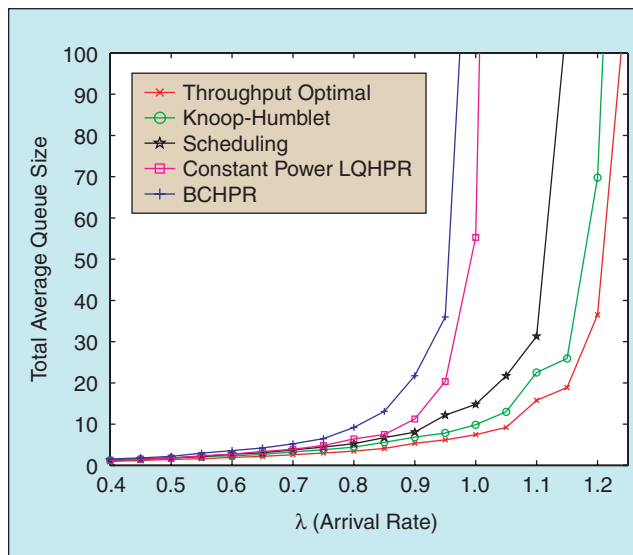
Assuming  $P^*(D)$  is strictly monotonic,  $D^*(P)$  will simply be its inverse. These quantities, like a rate-distortion curve in source coding, provide a characterization of the fundamental power delay tradeoff that can be achieved by any scheduling policy.

In [27], it is shown that  $P^*(D)$  is a strictly decreasing, convex function of  $D$ . From this it follows, that any point on this curve can be found by solving an average cost dynamic programming problem with a per stage cost given by  $J[n] = P(R[n], H[n]) + \beta (U[n]/\bar{I})$ , where  $\beta$  corresponds to a Lagrange multiplier for an average delay constraint. By varying  $\beta$ , different points on  $P^*(D)$  can be found. An example of  $P^*(D)$  is shown in Figure 4.

For a given  $\beta$ , the optimal rate allocation can be found numerically, and the optimal  $\beta$  for a given delay constraint can be found via standard convex programming techniques. General structural properties of the optimal policy can also be characterized [26], [27]. For example, assuming the fading is i.i.d., it can be shown that the optimal transmission rate is nondecreasing in the unfinished work. In [24], a version of this problem was considered where  $P(r, h)$  is linear in  $r$  for each  $h$ , as would be reasonable in the wide-band limit. This simplifies the dynamic programming problem and more precise structural results can be shown. In addition to considering average delay versus the long-term average power, there have been a number of papers which have examined related problems in a dynamic programming context e.g., [28]. In these papers, the object is again to minimize a cost related to energy, subject to various delay constraints, such as a deadline by which all packets must be sent.

### Asymptotic Power Delay Tradeoffs

The behavior of the power/delay tradeoff can be explicitly characterized in various asymptotic regimes. In [27] it is shown that as the average delay increases,  $P^*(D)$  approaches an asymptotic value of  $P_a(\bar{I})$  (see Figure 4). The limiting value  $P_a(\bar{I})$  is the minimum average power needed to support the average arrival rate, which corresponds to the average power level required for the channel to have a “long-term” capacity equal to  $\bar{I}/\Delta$ . Additionally, the rate at which  $P^*(D)$  converges to this asymptotic limit is given by  $1/D^2$  under any reasonable fading distribution. That  $P^*(D)$  can converge no faster than this is shown using drift arguments and the convexity of  $P(r, h)$ . This bound is shown to be tight by



▲ 4. An example of a power/delay tradeoff.

showing that a simple “threshold water-filling” policy [27] is order optimal, i.e., it converges to the asymptotic limit at the optimal rate of  $1/D^2$ . This policy has only a weak dependence on the buffer occupancy via a threshold rule. When the unfinished work exceeds this threshold, the transmitter uses a water-filling policy with an average rate that is greater than the arrival rate; when the unfinished work is less than the threshold, a water-filling policy with an average rate less than the arrival rate is used. The thresholds and the average rates used in each portion are adjusted depending on the average delay. It can also be shown that some dependence on the buffer occupancy is required for a policy to be order optimal [27]. In other words, the power required by any family of policies that do not depend on the buffer occupancy can not converge to  $P_n(\bar{I})$  at the rate of  $1/D^2$  as the average delay increases. This buffer dependence is needed for a policy to maintain a backlog of packets in the buffer. This backlog enables the transmitter to better exploit good channel conditions and smooths out some of the burstiness in the arrival process.

The behavior of  $P^*(D)$  can also be characterized in the regime of asymptotically small delays [29]. From the buffer dynamics in (6), the minimum possible average delay is one time-unit. By the small delay regime we mean the behavior of  $P^*(D)$  as  $D \rightarrow 1$ . The asymptotic value of  $P^*(D)$  in this regime depends on the fading distribution. Two distinct cases can be identified. The first case corresponds to channels such a Rayleigh fading channel where  $P^*(1)$  is infinite. This corresponds to a channel having a delay-limited capacity of zero [14]; i.e., using finite power it is impossible to send at a nonzero rate in every channel state. For these channels,  $D^*(P)$  will be greater than one for all finite  $P$ , but as  $P \rightarrow \infty$ ,  $D^*(P)$  will approach this value. The second case corresponds to those channels with  $P^*(1) < \infty$ . In this case,  $P^*(1)$  corresponds to the minimum average power needed for the channel to have a delay-limited capacity of  $\bar{I}/\Delta$ . Once again, the rate at which  $D^*(P)$  approaches its asymptotic limit can be characterized [29]. This also depends on the fading distribution, in particular on its behavior near zero. For example, in the case of Rayleigh fading,  $D^*(P) - 1$  decreases at a rate of  $e^{-\alpha P}$ , where  $\alpha$  is a constant depending on the fading distribution near zero. This is much faster than the optimal rate in the large delay regime, which implies that the reduction in power by a slight increase in tolerable delay is very significant in the small delay regime. An order optimal policy in this regime is a “channel-threshold policy” which transmits at a fixed rate, whenever the channel gain is greater than a threshold. Interestingly, this policy requires no dependence on the buffer state.

These results are robust to many variations in the model. For example, the power cost in (7), may be replaced by any convex increasing function; this can be used to reflect the power required for specific modulation and coding schemes. The results can be extended to models that allow the possibility of pack-

et errors and retransmissions as well as models with finite buffer sizes [30].

### **Multiaccess Model**

The above approach can be generalized to a multiuser setting, where a centralized control policy specifies the transmission rate and power for each user [30]. For example, consider the multiaccess channel with  $M$  users as in (1). In this case, let  $\mathcal{R}(\mathbf{h}, \mathbf{u})$  denote a rate allocation policy that specifies the transmission rate of each user as a function of the joint fading states and buffer occupancies. Now, the corresponding power cost will be given by

$$P(\mathbf{r}, \mathbf{h}) = \inf \left\{ \sum_{i=1}^M w_i P_i \text{ such that } \mathbf{r} \in \mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathbf{p}) \right\}, \quad (10)$$

where  $w_i$  are given weights, and  $\mathcal{C}_{\text{MAC}}(\mathbf{h}, \mathbf{p})$  denotes the multiaccess capacity region in (2) corresponding to power allocation  $\mathbf{p} = (p_1, \dots, p_M)$ . A solution to the optimization can always be found such that the rate vector  $\mathbf{r}$  lies at one of the  $M!$  extreme points of the capacity region.

In this case, a natural generalization of the optimal power/delay tradeoff is to define  $P^*(D_{\text{sum}})$  to be the minimum average weighted sum power required for the (possibly weighted) sum of the queuing delays to be no greater than  $D_{\text{sum}}$ . Once again, the power delay tradeoff can be evaluated via a dynamic programming formulation. The asymptotic analysis also carries through to this case. For example, as the average weighted sum delay,  $D_{\text{sum}}$  increases,  $P^*(D_{\text{sum}})$  can be shown to decrease to its asymptotic limit at a rate of  $1/D_{\text{sum}}^2$ . An order optimal sequence of policies can again be specified that require only a weak dependence on each user’s queue state. Specifically, a centralized controller only requires one bit of information about each user’s queue to implement these policies. This information will again indicate whether the queue size is above or below a threshold. Given this information the controller can identify one of  $2^M$  quadrants within which the joint buffer state lies. The controller then implements a policy for each quadrant that depends only on the channel state  $\mathbf{h}$ . This provides insight into the amount of control information that must be shared among users in a distributed setting to implement an order optimal policy.

### **Offline and Look-Ahead Scheduling Algorithms**

We briefly mention another approach for energy efficient scheduling from [25]. A finite horizon problem is considered with a deadline of  $T$  seconds. During time  $[0, T)$  packets randomly arrive and all packets that arrive in this interval must be transmitted over a channel without fading by time  $T$ . The goal is to accomplish this using the minimum energy. This is done by specifying the transmission time per packet,  $\tau$ . For a given transmission time  $\tau$ , the required energy is given by  $\varepsilon(\tau)$ ,

where  $e(\tau)$  is a decreasing convex function. Again,  $e(\tau)$  can be related to the channel capacity. In [25], an optimal “offline” scheduling algorithm is first considered, where all packet arrivals times are known a priori. This results in a convex optimization problem, which admits a simple solution. Approximate “online” algorithms based on a look-ahead buffer are developed which exploit the structure of the offline algorithm. These online algorithms are shown to require an average energy quite close to the optimal offline algorithm via simulations. This approach has the advantage that it does not require detailed knowledge of the arrival statistics.

## Acknowledgments

This research was supported in part by NSF Grants CCR-0313329 and CCR-0313183 and by ARO Grant DAAD19-03-1-0229.

*Randall A. Berry* received the B.S. degree in electrical engineering from the University of Missouri-Rolla in 1993 and the M.S. and Ph.D. degrees in electrical engineering and computer science from the Massachusetts Institute of Technology (MIT) in 1996 and 2000, respectively. He is currently an assistant professor in the Department of Electrical and Computer Engineering at Northwestern University. In 1998 he was on the technical staff at MIT Lincoln Laboratory in the Advanced Networks Group. His primary research interests include wireless communication, data networks, and information theory. He is the recipient of a 2003 NSF CAREER award.

*Edmund M. Yeh* received his B.S. in electrical engineering with distinction from Stanford University in 1994. In the same year, he received a Winston Churchill Scholarship. He received his master of philosophy in electrical engineering from the University of Cambridge in 1995 and his Ph.D. in electrical engineering and computer science from MIT in 2001. Since July 2001, he has been an assistant professor of electrical engineering and computer science at Yale University, New Haven, Connecticut. He received the Army Research Office Young Investigator Program Award and the National Science Foundation and Office of Naval Research Fellowships for graduate study. He is a member of Phi Beta Kappa, Tau Beta Pi, and IEEE.

## References

- [1] R. Gallager, “A perspective on multiaccess channels,” *IEEE Trans. Inform. Theory*, vol. 31, no. 2, pp. 124–142, 1985.
- [2] I.E. Telatar and R. Gallager, “Combining queuing theory with information theory for multiaccess,” *IEEE J. Selected Areas Commun.*, vol. 13, no. 6, pp. 963–969, 1995.
- [3] D. Bertsekas and R. Gallager, *Data Networks*. NJ: Prentice Hall, 1992.
- [4] Q. Zhao and L. Tong, “A multi-queue service room MAC protocol for wireless networks with multipacket reception,” *IEEE/ACM Trans. Networking*, vol. 11, Feb. 2003, pp. 125–137.
- [5] E. Yeh, “An inter-layer view of multiaccess communications,” in *Proc. 2002 ISIT*, Lausanne, Switzerland, 2002, p. 112.
- [6] E. Yeh and A. Cohen, “Throughput and delay optimal resource allocation in multiaccess fading channels,” in *Proc. Int. Symp. Information Theory*, Yokohama, Japan, 2003, p. 245.
- [7] E. Yeh and A. Cohen, “Information theory, queuing, and resource allocation in multi-user fading communications,” in *Proc. Conf. Information Sciences and Systems*, Princeton, NJ, 2004, pp. 1396–1401.
- [8] E. Yeh, “Delay optimal multiaccess communication for general packet length distributions,” in *Proc. Int. Symp. Information Theory*, Chicago, IL, 2004, p. 247.
- [9] A. Goldsmith and P. Varaiya, “Capacity of fading channels with channel side information,” *IEEE Trans. Inform. Theory*, vol. 43, no. 6, pp. 1986–1992, Nov. 1997.
- [10] R. Knopp and P. Humblet, “Information capacity and power control in single-cell multiuser communications,” in *Proc. Int. Conf. Communications*, Seattle, WA, 1995, pp. 331–335.
- [11] T. Cover and J. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
- [12] R. Laroia, S. Uppala, and J. Li, “Designing a broadband mobile wireless network,” *IEEE Signal Processing Mag.*, vol. 21, no. 5, pp. 20–28, 2004.
- [13] M.J. Neely, E. Modiano, and C.E. Rohrs, “Power and server allocation in a multi-beam satellite with time varying channels,” in *Proc. Infocom 2002*, New York City, 2002, pp. 138–152.
- [14] D. Tse and S. Hanly, “Multi-access fading channels: Part I and Part II,” *IEEE Trans. Inform. Theory*, vol. 44, no. 7, pp. 2796–2831, 1998.
- [15] L. Tassiulas and A. Ephremides, “Dynamic server allocation to parallel queues with randomly varying connectivity,” *IEEE Trans. Inform. Theory*, vol. 39, no. 2, pp. 466–478, 1993.
- [16] M. Andrews, K. Kumaran, K. Ramanan, A. Stolyar, and P. Whiting, “Providing quality of service over a shared wireless link,” *IEEE Commun. Mag.*, vol. 39, no. 2, pp. 150–154, 2001.
- [17] R. Leelahakriengkrai and R. Agrawal, “Scheduling in multimedia CDMA wireless networks,” *IEEE Trans. Vehicular Technol.*, vol. 52, no. 1, pp. 226–239, 2003.
- [18] S. Shakkottai and A.L. Stolyar, “Scheduling for multiple flows sharing a time-varying channel: The exponential rule,” *Trans. Aer. Math Soc. Transl. Ser. 2*, vol. 207, pp. 185–201, 2002.
- [19] A. Eryilmaz, R. Srikant, and J. Perkins, “Stable scheduling policies for broadcast channels,” in *Proc. ISIT*, 2002, p. 382.
- [20] A. Marshall and I. Olkin, *Inequalities: Theory of Majorization and Its Applications*. New York: Academic, 1979.
- [21] D. Tse, “Optimal power allocation over parallel gaussian broadcast channels,” in *Proc. 1997 Int. Symp. Information Theory*, Ulm, Germany, 1997, p. 27.
- [22] L. Li and A. Goldsmith, “Capacity and optimal resource allocation for fading broadcast channels: Part I: Ergodic capacity,” *IEEE Trans. Inform. Theory*, vol. 47, no. 3, pp. 1083–1102, 2001.
- [23] S. Verdú, “On channel capacity per unit cost,” *IEEE Trans. Inform. Theory*, vol. 32, pp. 1019–1030, 1990.
- [24] B. Collins and R. Cruz, “Transmission policies for time varying channels with average delay constraints,” in *Proc. 1999 Allerton Conf. Commun., Control, & Comp.*, Monticello, IL, 1999.
- [25] B. Prabhakar, E. Uysal-Biyikoglu, and A.E. Gamal, “Energy-efficient transmission over a wireless link via lazy packet scheduling,” in *Proc. IEEE Infocom 2001*, 2001, pp. 386–394.
- [26] V. Goyal, A. Kumar, and V. Sharma, “Power constrained and delay optimal policies for scheduling transmission over a fading channel,” in *Proc. IEEE Infocom 2003*, San Francisco, Mar. 30–Apr. 3, 2003, pp. 311–320.
- [27] R. Berry and R. Gallager, “Communication over fading channels with delay constraints,” *IEEE Trans. Inform. Theory*, vol. 48, pp. 1135–1149, May 2002.
- [28] A. Fu, E. Modiano, and J. Tsitsiklis, “Optimal energy allocation and admission control for communications satellites,” *IEEE/ACM Trans. Networking*, vol. 11, pp. 488–501, June 2003.
- [29] R. Berry, “Optimal power/delay tradeoffs in wireless communications—small delay asymptotics,” in *Proc. IASTED CIIT*, Scottsdale, AZ, Nov. 2003.
- [30] R. Berry, “Communication over fading channels with finite buffer constraints—Single user and multiple access cases,” in *Proc. 2001 IEEE Int. Symp. Info. Theory*, Washington, DC, June 24–29 2001, p. 58.