Dynamic Bandwidth Allocation Algorithms for High-Speed Data Wireless Networks

Matthew Andrews, Simon C. Borst, Francis Dominique, Predrag R. Jelenkovic, Krishnan Kumaran, K. G. Ramakrishnan, and Philip A. Whiting

Next-generation wireless networks are expected to support a wide range of highspeed data services, with Web browsing as one of the major applications. Although high data rates have been shown feasible in a single-user setting, the resource allocation issues that arise in a multiple-user context remain extremely challenging. Compared with voice, data traffic is typically more bursty, while the users are less sensitive to delay. These characteristics require resource allocation strategies to operate in a fundamentally different manner if the spectrum is to be used efficiently. In this paper we propose several algorithms for scheduling the efficient transmission of data to multiple users. As a new feature, the various schemes exploit knowledge of the buffer contents to achieve high throughput, while maintaining fairness by providing quality of service (QoS) to individual users. The proposed algorithms are backward compatible with existing cellular and personal communications services (PCS) standards such as IS-136. They provide a powerful approach to improving spectrum efficiency in forthcoming high-speed data cellular services. The extensive simulation experiments we present in this paper demonstrate that the algorithms significantly outperform conventional schemes.

Introduction

Present mobile cellular systems have only a limited ability to transmit data, a fact that has been observed by I. Alanko et al.¹ and C. J. Mathias.² For example, the Global System for Mobile Communications (GSM) and IS-95 standards permit a transmission rate up to about 10 kb/s. This rate, however, is too low to satisfy the requirements of many applications. In this paper we present algorithms for dynamically allocating carriers to mobile units to provide peak rates on the order of 400 kb/s. These algorithms can be implemented in conjunction with a group of narrowband time division multiple access (TDMA) carriers, such as the North American IS-136 system or the European GSM system. Indeed, they can be used in existing networks and operated within modest bandwidth requirements of a few megahertz.

To understand our approach, the reader must bear in mind two important distinctions between voice and data traffic. First, data traffic is much more bursty than voice traffic. Second, data traffic has different qualityof-service (QoS) requirements than those of voice traffic. On the one hand, data traffic is much more tolerant to delay than voice traffic; response times of several seconds are acceptable for data, whereas voice signals can only be delayed a fraction of a second. On the other hand, data is much more sensitive to bit errors, requiring error rates of 10⁻⁶, whereas voice can be satisfactorily transmitted with 10⁻³.

These distinctions led us to very different approaches for allocating bandwidth. In particular, the

allocation of a dedicated connection to a data user is less reasonable than to a voice user. A low-speed connection produces delays that are unacceptably high; a high-speed connection reduces delay, but cannot be used efficiently because of the bursty nature of data traffic. However, the tolerance of data applications to delay makes it possible to coordinate base stations and to schedule transmissions. Such an approach has also been suggested by J.-P. M. G. Linnartz,³ who proposed tight synchronization at the time slot level, to be used along with a reservation scheme. Linnartz also describes an alternative collision resolution approach.³

Of course, sophisticated scheduling would not be necessary if fixed channel assignment (FCA) provided an adequate QoS at a cost acceptable to network operators. This is not the case, however. In FCA, the base stations are partitioned into *reuse groups*—that is, groups of mutually non-interfering cells—and the carriers are statically divided among the reuse groups. Because the carriers are allocated on a fixed basis, the approach lacks flexibility. When the number of mobile units has reached the maximum that can be supported by the carriers allocated to a base station, further requests for service must be rejected. In reality, of course, a simple reallocation of carriers might have made it possible to support these extra mobile units at little cost to other network users.

A similar difficulty arises if packets build up at a particular base station; this can occur even if the mobile units are evenly distributed. Once again, extra carriers cannot be allocated because the number at each base station is fixed, even though few packets are being presented for transmission at base stations nearby. This is in contrast to voice systems, where the peak bandwidth requirements are fixed at a relatively low level. Thus, the inefficiency of FCA in allocating spectrum precludes its use by operators.

The algorithms we present here are for wireless networks operated according to a frequency reuse plan, but they do allow for dynamic channel assignment (DCA). In DCA, carriers are not statically assigned to base stations; they may be diverted in a dynamic fashion from cells with traffic lulls to cells with traffic peaks. A frequency reuse plan limits the interference between co-channel users by geographic

Panel 1. Abbreviations, Acronyms, and Terms
DCA—dynamic channel assignment FCA—fixed channel assignment GSM—Global System for Mobile Communications ILP—integer linear program LP—linear program MSC—mobile switching center PCS—personal communications services QoS—quality of service TDMA—time division multiple access

distance according to worst-case conditions, making it unnecessary to obtain on-line propagation or interference measurements. The sole problem remaining is scheduling, which is subject to these reuse constraints. Unfortunately, the optimal solution to such a scheduling problem is not known. Furthermore, even if the solution were known, its complexity would make it impractical to implement.

In our approach to this problem, we devised a number of heuristics, each of which determines carrier allocations to base stations for a fixed period based on the number of packets awaiting transmission. Once the carriers are allocated, a common algorithm determines how they are used to transmit the mobile packets. J. M. Harris and S. P. Kumar⁴ address the related and easier question of how to schedule the transmission of a given set of packets onto carriers whose allocation to base stations has already been determined. The model in Harris and Kumar⁴ includes the possibility that mobile units are able to receive packets from more than one base station, which is not considered here. Closely coupled in the circumstances described in Harris and Kumar are the dual problems of how to assign carriers to base stations and then assign these carriers to mobile units.

The rest of this paper is organized as follows. First we describe the model set-up and the issues surrounding our choice for the source model. The next three sections propose a number of heuristic carrier allocation algorithms defined within a common format. In the section that follows, we evaluate the communication requirements and computational complexity of the various algorithms. Next, we present a series of simulation experiments, in which we consider the impact of the burstiness and the duration of the scheduling interval on the delay and throughput performance. We conclude with our expectations for next-generation wireless networks.

Model Description

We consider a wireless network of N base stations supporting M users. Packets destined for the users arrive at the serving base station, where they are queued until they are transmitted. (Throughout, we focus on the stream of packets from the base station down to the users. We do not consider the flow of packets from the users up to the base station.) The base stations share a pool of K orthogonal carriers for transmitting packets down to the users. Transmission occurs in a slotted fashion, with exactly one packet transmitted on a single carrier during a time slot.

Base stations may reuse carriers subject to certain interference constraints. We assume the reuse constraints may be described in the form of an *interference graph*, with the vertices corresponding to the base stations and the edges representing the pairs of interfering base stations. Thus, if the interference graph contains the edge { j_1 , j_2 }, then base stations j_1 and j_2 are barred from transmitting on the same carrier in the same time slot. Such models of networks of queues, whose service constraints are determined by a graph, also appear in the context of input-queued switches, as described by L. Tassiulas.⁵

The scheduling policy determines the assignment of carriers to users in each time slot. Let $X_{jkt} \in \{0,1\}$ indicate whether carrier *k* is allocated to base station *j* in the *t*-th time slot. For the assignment to be feasible, the reuse constraints require $X_{j_1kt} + X_{j_2kt} \leq 1$ for all edges $\{j_1, j_2\}$ in the interference graph.

Let $Y_{ikt} \in \{0,1\}$ indicate whether carrier *k* is assigned to user *i* in the *t*-th time slot. Then $\sum_{i \in M_j} Y_{ikt} \leq X_{jkt}$ for all j = 1, ..., N, where M_j represents the set of all users served by base station *j*. Define $Z_{it} := \sum_{k=1}^{K} Y_{ikt}$ as the number of carriers assigned to user *i* in the *t*-th time slot. A user may not be assigned more than *L* carriers simultaneously, that is, $Z_{it} \leq L$.

Denote by Q_{it} the queue length for user *i* at the

start of the *t*-th time slot, that is, the number of packets for user *i* waiting for transmission. The queue lengths evolve over time as

$$Q_{i(t+1)} = A_{it} + [Q_{it} - Z_{it}]^+$$

where $[z]^+ := \max \{z, 0\}$ and A_{it} is the number of packets that arrive for user *i* during the *t*-th time slot.

Suppose that a packet for user *i* arrives in a queue of length $Q_{it'}$ during the *t'*-th time slot. The delay experienced by that packet, measured in time slots, is then defined to be (t'' - t'), with

$$t'' := \min\left\{t : \sum_{t=t'+1}^{t} Z_{it} \ge Q_{it'}\right\}.$$

We now discuss the traffic model. The most basic statistical model that can capture data "burstiness" and a complex dependency structure is the so-called "on-off model," first investigated by J. W. Cohen⁶ and M. Rubinovitch.⁷ For two-state Markov (fluid) on-off models, D. Anick, D. Mitra, and M. M. Sondhi⁸ analyzed the impact of the burstiness on queuing performance. Subsequent studies explored more general Markov models with finite state space. These led to the equivalent-bandwidth theory for Markovian (or in general exponentially bounded) arrival processes, which was treated by A. Elwalid et al.⁹ and P. W. Glynn and W. Whitt.¹⁰

Recently, statistical analysis has provided increasing evidence that the traffic streams in modern broadband networks exhibit long-tailed (subexponential) characteristics. The types of statistical results generated by W. Willinger et al.¹¹ for Ethernet traffic have stimulated research in queuing analysis under the heavytailed (non-Cramér) assumptions. P. R. Jelenkovic and A. A. Lazar¹² have recently published their results on multiplexing on-off sources with subexponential characteristics, along with an extensive list of references on subexponential models.

It is difficult to predict the predominant applications and statistical characteristics that will govern future networks. Consequently, we have conducted our experiments using both exponential and subexponential (long-range dependent) models.

First, we present a formal definition of an on-off process. Consider two independent sequences of i.i.d. ran-



Figure 1. Sample path of an on-off model.

dom variables $\{\tau_n^{off}, n \ge 0\}, \{\tau_n^{on}, n \ge 0\}, \tau_0^{off} = \tau_0^{on} = 0.$ Define a point process $T_n^{off} := \sum_{i=0}^n (\tau_i^{off} + \tau_i^{on}), n \ge 0.$ This process will be interpreted as representing the beginnings of off periods in an on-off process. Next, let $\{B_{\nu}, t = 0, 1, ...\}$ be a sequence of i.i.d. random variables.

Then, a discrete time on-off process A_t is defined as

 $A_t = B_t \text{ if } T_n^{off} - \tau_n^{on} \leq t < T_n^{off} , n \geq 1,$

and $A_t = 0$, otherwise. **Figure 1** presents a sample path for an on-off model.

Table I combines the various distribution functions that we used to model τ^{off} , τ^{on} , and B_t . We chose only the exponential distribution for the off periods because large queue build-ups are insensitive to the distribution of off periods (Jelenkovic and Lazar¹² give a rigorous treatment of this insensitivity phenomena). In contrast, the queuing behavior is strongly affected by the duration of the on times, with a clear dichotomy between exponential and subexponential distributions. Fluctuations in the peak data rate are captured with a Poisson distribution. In general, for all different choices of distributions and their parameters, we will show that DCA schemes significantly outperform the FCA approach.

Preliminaries

In view of the cell-by-cell reuse constraints, the problem of scheduling the transmission of packets to

Table I. Combinations of distributions used in Experiments 1 through 5.

$_{\mathcal{T}}$ off	Ton	B _t
Exponential	Exponential	Poisson
Exponential	Pareto	Poisson

users may be decomposed into two parts:

- Allocation of carriers to base stations; and
- Assignment of carriers to individual users, given the number of carriers allocated to the base stations.

In this paper, we focus on the first problem; the next two sections describe various schemes, all of which differ in the way they allocate carriers to base stations.

The approach to the second problem is common to all these schemes; we assume that the assignment of carriers to users is proportional to the queue lengths, subject to the constraint that no user may be assigned more than L carriers. Specifically, suppose base station j has been allocated K_j carriers. First, we assign

$$L_{i} = \min\left\{L, \left\lfloor K_{j}Q_{i} / \sum_{m \in M_{j}}Q_{m}\right\rfloor\right\}$$

carriers to each user $i \in M_j$ served by base station j. In case any carriers remain—that is, $\sum_{i \in M_j} L_i < K_j$ —we assign them to the user $i \in M_j$ with the largest ratio of Q_i / L_i among those with $L_i < L$. We repeat the above procedure until either no carriers remain, or $L_i = L$ for all users $i \in M_j$.

There are numerous ways to approach assigning

carriers to users. Since no user may be assigned more than L carriers, capacity may potentially be wasted if there are only a few users with non-empty queues. From the point of view of throughput, it seems attractive to assign carriers to the user with the longest queue. To be fair, however, it appears more reasonable to assign a comparable number of carriers to each user. The proportional assignment scheme described above may be viewed as an intermediate strategy with regard to these two extremes. Among the class of work-conserving schemes, we do not expect any significant variation in average-delay performance as long as the value of L is not extremely small.

Since the above procedure for assigning carriers to users does not require any exchange of information or complex calculations, we assume that the assignment is performed in every time slot. In contrast, we assume that carriers are allocated to base stations only once every *S* time slots, for some parameter *S*, because the process may involve communication overhead and extensive computations. We refer to such a scheduling interval of *S* time slots as a *superframe*. The parameter *S* may be different for different schemes, depending on the communication and computational complexity involved. A large value of *S* allows for a more elaborate allocation process, whereas a small value of *S* enables a more rapid response to congestion conditions.

Because each scheme uses a superframe structure, the base stations must be globally synchronized. This need not be very accurate if superframes are comparatively long—that is, if *S* is large. Of course, long superframes introduce some additional latency, but the increase in delay is not significant as long as the superframe is shorter than the mean burst period.

We now discuss average delay bounds. Determining the delay-minimizing allocation scheme is prohibitively demanding for all but the simplest cases. Hence, a bound for the achievable average delay is instrumental for evaluating the performance of heuristics.

A simple but effective delay bound may be derived as follows. Consider a clique of base stations in the interference graph—that is, a group of base stations that interfere with each other. By definition, carriers cannot be reused in the clique. Thus, the aggregate queue length in the clique cannot be any smaller than it would be, in case the aggregate packet arrival stream were served by K dedicated carriers in isolation. Hence, according to Little's law,¹³ the average packet delay in the clique is also bounded below by the average delay if the aggregate packet arrival stream were served by K isolated carriers.

Now let us suppose the interference graph is symmetric and the maximum clique size is *F*. In that case, the network-average delay is bounded below by the average delay if the superposition of *F* arrival streams were served by *K* isolated carriers. We refer to the latter delay bound as the *pooling bound*. Tighter bounds may be derived along similar lines.

Distributed Carrier Assignment

In this section we outline two distributed assignment schemes whose only requirement for allocating carriers is the local exchange of information between base stations. In the next section, we present a centralized allocation scheme, which uses global knowledge to assign carriers to base stations.

Both schemes use the concept of reuse groups. We assume that the base stations are partitioned into F reuse groups N_1, \ldots, N_F . A reuse group is an independent set in the interference graph—that is, a group of mutually non-interfering base stations.

In the first scheme, the carriers are also partitioned into *F* sets, $K_1, ..., K_{F_1}$ each associated with a particular reuse group. The carriers are designated, but not restricted, for use by the base stations in the corresponding reuse group. On a request basis, base stations may borrow carriers that are not currently needed by the designated owners. Henceforth we refer to this scheme as *distributed carrier requesting*.

In the second scheme, the carriers are not associated with any particular group of base stations; instead, they float freely. At the start of every superframe, the base stations in a particular reuse group grab all the carriers not currently being used by any interfering base stations. We refer to this scheme as *distributed carrier raking*.

Of course, it is conceivable to have a hybrid scheme in which some carriers are designated, while others are completely floating. For ease of presentation, however, we restrict our discussion to the two extreme schemes mentioned above.

Both schemes also use the concept of base station "demand" for carriers. At the start of every super-frame, base station j calculates the demand for carriers in the next superframe as

$$D_j = \left\lfloor \sum_{i \in M_j} \min\left\{L, \frac{Q_i}{S}\right\} \right\rfloor$$

The *demand* is the number of carriers that base station *j* could acquire without risking any unused capacity in the next superframe.

The section below describes the two proposed schemes in detail.

Distributed Carrier Requesting

In the distributed carrier requesting scheme, base stations may borrow carriers not currently needed by the designated owners. This scheme, performed at the start of every superframe, consists of the following three steps:

- Restricted retainment of designated carriers,
- Restricted acquisition of nondesignated carriers, and
- Unrestricted acquisition of designated carriers.

Restricted retainment of designated carriers. Base station *j* acquires min $\{D_j, K/F\}$ designated carriers, with *K/F* representing the number of carriers associated with each reuse group.

Restricted acquisition of nondesignated carriers. Next, base stations may borrow carriers not already claimed by the designated owners. The process is organized in *F*–1 rounds. In the *n*-th round, the base stations in reuse group N_f concurrently acquire the carriers in the set K_{f+n} that have not yet been claimed by any interfering base stations in reuse groups N_{f+1}, \dots, N_{f+n} . However, base stations do not acquire more carriers than their computed demand. By definition, the base stations in a reuse group are all mutually non-interfering, guaranteeing that the assignment of carriers will remain feasible throughout. To enhance performance in response to congestion conditions, however, the order in which base stations are allowed to acquire nondesignated carriers could be periodically changed or dynamically adjusted.

		Reuse groups				
		<i>N</i> ₁	N ₂	N ₃	N ₄	
Carrier sets	<i>K</i> ₁	A D ₀ R	B ₃ D ₃	B ₂ D ₂	B ₁ D ₁	
	K ₂	B ₁ D ₁	A D ₀ R	B ₃ D ₃	B ₂ D ₂	
	<i>К</i> 3	B ₂ D ₂	B ₁ D ₁	A D ₀ R	B ₃ D ₃	
	K ₄	B ₃ D ₃	B ₂ D ₂	B ₁ D ₁	A D ₀ R	

Figure 2. Distributed carrier requesting process.

Unrestricted acquisition of designated carriers. Since base stations do not acquire more carriers than their computed demand, some carriers may remain unclaimed at the end of the previous stage. These are finally reclaimed by the designated base stations.

Figure 2 depicts the distributed carrier requesting process for a scenario with F = 4 reuse groups. The letter A corresponds to the first stage, the restricted retainment of carriers by the designated base stations. Then, in D_0 , the base stations declare which of their designated carriers they wish to retain. The symbol B_n represents the *n*-th borrowing round during the second phase, the restricted acquisition of carriers by nondesignated base stations. The base stations in reuse group N_f then decide which of the carriers in the set K_{f+n} they wish to borrow among those that have not been declared in use by any of the interfering base stations in reuse groups $N_{f+1}, ..., N_{f+n}$. In D_n , the base stations in reuse group N_f declare which of those carriers they have decided to borrow. The letter R finally indicates the third stage, the unrestricted acquisition of carriers by the designated base stations.

Distributed Carrier Raking

In this scheme, the base stations periodically rake all the carriers that are not currently used by any interfering base stations. In each superframe, only the base stations in one reuse group are active, while the base stations in all other reuse groups are passive. The active reuse group alternates in a cyclic manner—that is, the base stations in reuse group N_f are active in the (n F + f)-th superframe, n = 1, 2, ... The scheme, executed again at the start of every superframe, consists of three steps:

- Restricted retainment of carriers by passive base stations,
- Unrestricted acquisition of carriers by active base stations, and
- Unrestricted retainment of carriers by passive base stations.

Restricted retainment of carriers by passive base stations. First, the passive base stations decide which of the carriers they used in the previous superframe they wish to keep. Suppose that base station jused K_j carriers in the previous superframe. If $D_j < K_j$, then base station j retains D_j of these carriers and releases the remaining $K_j - D_j$ ones. If $D_j \ge K_j$, then base station j holds on to all K_j carriers.

Unrestricted acquisition of carriers by active base stations. Next, the active base stations concurrently grab all the carriers that are not already claimed by any interfering base stations. By definition, the base stations in a reuse group are all mutually non-interfering, guaranteeing that the assignment of carriers will remain feasible.

Unrestricted retainment of carriers by passive base stations. Finally, the base stations in the passive reuse group reclaim those carriers they relinquished in the first stage that were not claimed by any of the interfering active base stations.

Figure 3 illustrates the distributed carrier raking process for a scenario with F = 4 reuse groups. We assume that it is the turn of the base stations in the second reuse group to rake carriers. The letter A corresponds to the first stage, the restricted retainment of carriers used in the previous superframe by the passive base stations. Then, in D₀, the passive base stations declare which of the carriers used in the previous superframe they wish to retain. The letter C represents the second phase, the unrestricted acquisition of carriers by the active base stations. In D₁, the active base stations declare which of the carriers they have been able to acquire. The letter R indicates the third stage, the unrestricted



Figure 3. Distributed carrier raking process.

retainment of carriers used in the previous superframe by the passive base stations.

Centralized Carrier Allocation

In the previous section, we presented two distributed assignment schemes for allocating carriers, both of which involve only local exchange of information between base stations. We now outline a centralized allocation scheme, which relies on global knowledge for assigning carriers to base stations.

Successively, for each carrier, the scheme attempts to find a subset of mutually non-interfering base stations with the maximum number of packets queued. In reference to the associated interference graph, we term such a set as the *maximum-weight independent set*. The scheme is motivated by the stability results for the maximum-weight independent set established by N. Kahale and P. E. Wright.¹⁴ As before, carriers are allocated at the start of each superframe, and the capacity of the carrier is subtracted from the queues to which an assignment has been made.

The maximum-weight independent set problem can be formulated as an integer linear program (ILP), with variables $x_j \in \{0,1\}$, indicating whether base station *j* belongs to the desired set. In general, the problem is known to be NP-complete, even on very simple graphs, and we expect this to be the case on most wireless interference graphs as well. Therefore, we "relax" the integrality constraints and consider the corresponding linear program (LP). First we assign weights $\{w_j\}$ to each of the base stations j = 1, ..., N determined by $w_j = \sum_{i \in M_j} Q_i$ —that is, the aggregate queue length of the users served by base station *j*. Then, we successively perform the following steps for each of the carriers k = 1, ..., K.

1. Solve the linear program

$$\max \sum_{j=1}^{N} w_{j} x_{j}$$

sub
$$\sum_{l \in C} x_{l} \le 1 \text{ for all maximal cliques } C,$$

$$0 \le x_j \le 1$$
 for all $j = 1, \dots, N$.

- 2. Let $\{x_j^*\}$ be the solution to the linear program. Initialize $J := \{1, ..., N\}$ to the set of all base stations. Then perform the following steps.
 - (a) Let $\hat{j}:= \arg \max_{j \in J} \{x_j^*\}$. Assign carrier k to base station \hat{j} —that is, $X_{jk} := 1$. Remove base station \hat{j} and all its neighbors from the set J.

(b) Repeat the above step until $J = \emptyset$.

3. For each base station that has been assigned carrier k, reduce its weight by S—that is, w_i := w_j-SX_{jk}.

We can construct the maximal clique *C* in step 1 of the algorithm by repeatedly appending vertices to a set initialized to some vertex { *j* } so that each new vertex has edges to all the current elements of the set, until no more vertices can be added. Enumerating all maximal cliques is a formidable task on arbitrary graphs. However, the local nature of the interference in a wireless network ensures that the number of distinct maximal cliques stays within a constant factor of the number of cells in the network, thus rendering a linear program with a tractable number of constraints. Determining the true maximum-weight independent set requires us to constrain x_j to be binary, that is, 0 or 1. In the absence of these constraints, the LP solution.

Step 2 is essentially a method of "rounding" the



Figure 4. An odd cycle with seven base stations.

fractional solution of the LP to a feasible assignment of carriers. This produces a lower bound to the optimal value, which, in conjunction with the abovementioned upper bound, makes the LP a useful bounding procedure. In step 3, for each base station that has been assigned carrier k, we reduce the corresponding weight by the number of packets S that can be transmitted during the next superframe.

For well-structured interference graphs and random weights, the LP method provides a reasonable approximation to the true maximum-weight independent set. In experiments carried out with a regular four-cell reuse pattern, which are discussed in detail in a later section, the LP almost always produced the exact optimal solution. This was evident from the fact that the values of $\{x_j^*\}$ returned by the LP were integral, despite the relaxation.

However, there are cases in which the LP can be shown to "fail"—that is, to provide ambiguous values of x_j^* with no natural rounding. The simplest such cases are so-called "induced odd cycles," as shown in **Figure 4**. Consider a snapshot of the network where the only non-empty queues are located on such a ring with an odd number of base stations, and the queues are all equal in size. It can then be shown, by straightforward substitution, that the LP solution would be $x_j^* = 0.5$ for all *j*.

The true solution would have $x_j^* = 1$ on p base stations in a cycle of length 2p + 1, with $x_j^* = 0$ elsewhere. Such cases, although pathological, can arise in small subnetworks and must be dealt with by adding constraints to the LP of the form $\sum_{j \in C} x_j \leq p$. These constraints rule out the ambiguous solution, which has $\sum_{j \in C} x_j = p + 1/2$. Incorporating all such constraints on an arbitrary graph is prohibitive, but the effort is moderate for small values of *p* on wireless interference graphs. Even here we cannot cover all values of *p*, but this is not likely to be too detrimental in practice, since the effect of large values of *p* only occurs in the unlikely adversarial circumstance of a large ring of base stations containing the only non-empty queues.

Thus, the LP provides a heuristic for allocating carriers in a stable, efficient manner while avoiding interference, although the algorithm, as described above, may unfairly favor users or base stations with very long queues. The potential unfairness may be partly ameliorated by modifying the weights to reflect the desired level of fairness among users, but possibly at the expense of throughput. The ultimate solution may require integration of the scheduling algorithm with policing and admission control.

Complexity

The complexity of the various algorithms is determined by the level of communications required to pass messages between various network elements and the amount of computing effort needed to perform the calculations. Carrier raking and carrier requesting have a similar level of complexity, with the centralized algorithm being significantly more complex, both in communication and computation.

We first consider carrier requesting. A base station need only communicate with its neighbors. As we showed earlier in Figure 2, just before the start of a round, each base station broadcasts the current list of available carriers. Thus, in the first round, each base station declares which of the carriers still remain from those designated for its use. In the second round, each base station declares which carriers it is using from the group of carriers offered in the first round. Similar declarations are made in the third round, and so on.

The scheme may be implemented by using a common control channel, which would then be used for a total of F–1 rounds. If the reuse for signaling is H, then the signaling load for each round is HFK/F = HK bits, where K/F bits are needed to indicate which designated carriers are being used by each base station. The total number of bits that must be transmitted on the common control channel is *HFK*, since there are *F*–1 rounds and a declaration at the end. For example, if F = 4, H = 12, and K = 32, then 1,536 bits must be transmitted per superframe, so that if the capacity is 10 kb/s, the common control channel can process as many as 6 superframes per second. To avoid the need to send message and base station identifiers, a common control channel requires synchronized base stations. The computing resources required to calculate the demand and carry out the algorithm are very limited. Only a slight delay is incurred by the base stations' need to compute the number of carriers they will borrow before the next round is started.

In carrier raking, the passive base stations send out a broadcast signal once per superframe, indicating which carriers they are prepared to relinquish. Active base stations respond by declaring which carriers they have taken in a subsequent broadcast message. Once again, the signaling load is *HFK* bits per superframe if a common control channel is used. The computational load is similar to carrier requesting. The demand calculations also correspond to the ones used for carrier requesting, but these only have to be performed once every *F* superframes.

Last, we consider the centralized scheme. If the controller is located at the mobile switching center (MSC), then there will be no need to pass queuing information from the base stations, as far as the down link is concerned. (This would have to be done for the up link.) Once the succession of linear programs has been solved, the results must be passed to the base stations. For each station in the network, the results may be passed as a *K*-bit field and supplied with a base station identifier *B* bits long. Therefore, the total information that must be passed per superframe is (K+B)N, where *N* is the number of base stations. This information can be passed through wireline connections between the MSC and the base stations. No wireless capacity is used.

There are N variables in the linear program, and there are GN constraints, where G is a small factor that depends on the number of maximal cliques. The LP must be solved once per carrier, but the solution time can be shortened by using the previous solution as the

	Experiment 1	Experiment 2	Experiment 3	Experiment 4	Experiment 5
On period(s)	Varying	Varying	4	4	4
Off period(s)	Varying	10 * on period	36	36	36
On+off period(s)	40	Varying	40	40	40
Peak rate (kb/s)	Varying	384	384	384	384
Burst size (Mb)	1.536	Varying	1.536	1.536	1.536
Mobile units per base station	4	4	Varying	Non-uniform	4
Superframe(s)	0.24/1.00	0.24/1.00	0.24/1.00	0.24/1.00	Varying

Table II. Simulation parameters.

starting point for the next. There are several efficient methods for solving such LPs.

Simulation Experiments

We now present an overview of the results from our simulation experiments to examine the performance of the various algorithms proposed in the previous sections. We also present results for the pooling bound (described in the section titled "Preliminaries," earlier in this paper), which helps us set a lower bound on the achievable delay.

In all experiments, we consider a network of N = 64 base stations, arranged in a grid as depicted in **Figure 5**. We assume the grid is "wrap-around" (that is, a periodic boundary) to reduce the edge effects. The total number of carriers in the system is K = 32, and the maximum number of carriers that can be assigned to any particular mobile unit is L = 8. The bandwidth per carrier is 48 kb/s, and the packet size is 1,920 bits, making it possible, for example, to transmit S = 25 packets per second.

For the first four experiments we used a superframe of 1 sec. for the carrier requesting algorithm and the LP-based algorithm. For the carrier raking algorithm we used a superframe of 0.24 sec., so that each base station has the opportunity to acquire carriers about once every second.

In all experiments, we adopted the on-off traffic model described earlier. This model was used to describe the arrival of packets that need to be transmitted to a particular mobile unit. The key parameters in the traffic model are the average on period, the average off period, the mean rate, the peak rate, and the burst size. These quantities are related through



Figure 5. The interference graph.

$$Mean \ rate = \frac{On \ period \times peak \ rate}{On \ period + off \ period}$$

and

Burst size = $On period \times peak rate$.

To investigate the impact of the traffic characteristics on the performance of the various algorithms, we varied these parameters as summarized in **Table II**. The mean rate is always taken to be 38.4 kb/s. The number of packet arrivals per slot is assumed to have a Poisson distribution, with the peak rate divided by the packet size as the parameter. As will be seen, all the experiments show that in general the DCA schemes are supe-



Figure 6. Packet delay with fixed burst size (exponential on period).



Figure 7. Packet delay with fixed burst size (Pareto on period).

rior, no matter whether we take the mean delay or the tail distribution of the delay (95th percentile) as the performance criterion. This is true even when the distribution of mobile units is uniform, the case most favorable to FCA. Below we describe the results from the simulation experiments in greater detail.

Experiment 1. Fixed Burst Size for Exponential/Pareto On Periods

For this experiment, we had four mobile units per base station. We constructed the traffic model so that the "bursts of packets" would have a fixed mean size. However, we altered the peak arrival rate so that these bursts would have different shapes.

Figures 6 and **7** show the average delay as a function of the on period. (Note that an increasing on period corresponds to a decreasing peak rate and decreasing burstiness.) As we expected, the average delay decreases as the burstiness decreases. Furthermore, the schemes are more set apart when the traffic is more bursty. This phenomenon illustrates the responsiveness of the dynamic schemes. However, when the burstiness is reduced, both carrier raking and LP exhibit a significant residual latency, caused by the continuing response to queue fluctuations, despite the near-constant traffic.

Experiment 2. Fixed Peak Rate in Exponential On Periods

For this experiment, we also had four mobile units per base station. However, the peak rate was fixed during the on period, whose length was altered to vary the burst size. Because the peak rate and load are fixed, burstiness is not changing in this case.

As **Figure 8** shows, the delay worsens somewhat with the mean length of the on period. This reflects the fact that more packets are queued up in a burst as its duration increases.

Experiment 3. Varying Number of Mobile Units per Base Station for Exponential/Pareto On Periods

We kept the burst size and peak rate fixed, and varied the number of mobile units per base station. This has the effect of varying the total load on the system.

Figure 9 indicates the gains made in throughput from using DCA. If we assume a mean delay requirement of a few seconds, for example, then we can

attain throughput improvements of around 30 to 70%. Carrier requesting and LP have similar throughput, with carrier raking being a little worse. The gains in throughput over FCA are greater when the traffic is more bursty, as shown in **Figure 10**.

Experiment 4. Non-Uniform Distribution of Mobile Units for Exponential On Periods

This experiment was the same as experiment 3, except that we had a non-uniform distribution of mobile units. If the average number of mobile units per base station was *m*, then we assigned to each base station a number of mobile units that was chosen with equal probability from $\{m-2, ..., m+2\}$.

Figure 11 shows the impact of this non-uniformity on performance. Here FCA delay performance worsens significantly as the mismatch between the number of carriers allocated and the actual number of mobile units at a base station increases. Carrier raking also does badly, because of the possibility of getting "locked" into a poor match between carriers and load, prompting base stations to hold on to their carriers.

Experiment 5. Changes in Superframe Duration for Exponential On Periods

For the DCA schemes, we altered the length of the superframes to observe the effect that latency in the carrier allocation process has on performance.

The results depicted in Figure 12 meet our expectations by showing that, for all the schemes, mean delay improves as the superframe duration drops. As we can see, a superframe duration on the order of a fraction of one second is very desirable. The complexity calculations in the previous section indicate that this should be achievable for both carrier raking and carrier requesting. The performance of carrier requesting degrades gradually as the duration of the superframe increases. The average delay approaches that of FCA as the superframe duration grows large, because the queue length becomes negligibly small compared to the superframe capacity, reducing the computed demand to zero. In contrast, the performance of carrier raking degrades severely as the superframe duration increases, causing excessive delay for superframes longer than one second.



Figure 8. Packet delay with fixed peak rate (exponential on period).



Figure 9. Packet delay with varying number of mobiles (exponential on period).



Figure 10. Packet delay with varying number of mobiles (Pareto on period).



Figure 11. Packet delay with varying (non-uniform) number of mobiles (exponential on period).





Conclusions

Next-generation wireless networks are expected to support a wide range of high-speed data services, with Web browsing being one of the major applications. Compared with voice, data traffic is typically more bursty, while the users are less sensitive to delay. These characteristics require resource allocation strategies to operate in a different manner in order to use spectrum efficiently. In particular, the allocation of dedicated bandwidth to data applications is less reasonable than to a voice user. If a lowbandwidth connection is provided, then the burstiness will cause excessive delay and loss of packets. If a high-bandwidth connection is established, then the delay and loss will be less serious, but the utilization will be poor. The delay tolerance of data applications, however, allows for the possibility of coordinating packet transmissions among base stations. Thus, the efficient operation of high-speed data wireless networks requires the use of dynamic bandwidth allocation algorithms.

We have proposed several such algorithms for coordinating scheduling of packet transmissions among base stations. As a new feature, the various schemes exploit knowledge of the buffer contents and achieve high throughput, while maintaining fairness by providing QoS to individual users. The proposed algorithms backward compatible with existing cellular and PCS standards such as IS-136—provide a powerful approach to improving spectrum efficiency in forthcoming high-speed data cellular services.

We have conducted extensive simulation experiments to demonstrate the efficiency of the algorithms. Not surprisingly, there is a tradeoff between the throughput and delay performance of the algorithms and the communication and computation overhead involved. Since in practice the available resources for these tasks are limited, the theoretically optimal strategies are not necessarily the most adequate algorithms from a practical perspective.

The algorithms that we considered adopt the concept of a frequency reuse plan. Strategies that use signal strength measurements to determine the positions of individual mobile units may achieve an even higher spectrum efficiency at the expense of larger operational complexity, as observed by J. C. Chuang and N. Sollenberger.¹⁵

References

- I. Alanko, M. Kojo, H. Laamanen, M. Liljeberg, M. Moilanen, and E. Raatikainen, "Measured performance of data transmission over cellular telephone networks," *Comp. Commun. Rev.* (USA), Vol. 24, No. 5, 1994, pp. 24–44.
- C. J. Mathias, "Wireless data—What's real? What's wrong?," *Bus. Commun. Rev.* (USA), Vol. 27, No. 6, June 1997, pp. 52–54.
- 3. J.-P. M. G. Linnartz, "On the performance of packet-switched cellular networks for wireless data communications," *Wireless Networks*, Vol. 1, No. 2, 1995, pp. 129–138.
- J. M. Harris and S. P. Kumar, "An algorithm for minimizing queueing delay of packets in cellular data networks," *Proc. 33rd Annual Allerton Conf. on Commun., Control, and Computing,* Monticello, Ill., Oct. 4–6, 1995, pp. 945–954.
- L. Tassiulas, "Linear complexity algorithms for maximum throughput in radio networks and input queued switches," *Proc. IEEE INFOCOM* '98, San Francisco, Apr. 1998, pp. 533–539.
- J. W. Cohen, "Superimposed renewal processes and storage with gradual input," *Stochastic Processes and Their Applications*, Vol. 2, No. 1, Jan. 1974, pp. 31–57.
- M. Rubinovitch, "The output of a buffered data communication system," *Stochastic Processes and Their Applications*, Vol. 1, No. 4, Oct. 1973, pp. 375–382.
- D. Anick, D. Mitra, and M. M. Sondhi, "Stochastic theory of a data handling system with multiple sources," *Bell Sys. Tech. J.*, Vol. 61, No. 8, Oct. 1982, pp. 1871–1894.
- A. Elwalid, D. Heyman, T. V. Lakshman, D. Mitra, and A. Weiss, "Fundamental bounds and approximations for ATM multiplexers with applications to video teleconferencing," *IEEE J. on Sel. Areas in Commun.*, Vol. 13, No. 6, Aug. 1995, pp. 1004–1016.
- P. W. Glynn and W. Whitt, "Logarithmic asymptotics for steady-state tail probabilities in a single-server queue," *J. of Applied Probability*, Vol. 31A, 1994, pp. 131–156.
- W. E. Leland, W. Willinger, M. S. Taqqu, and D. V. Wilson, "On the self-similar nature of Ethernet traffic," *Proc. SIGCOMM '93*, San Francisco, Sept. 13–17, 1993, pp. 183–193.
- 12. P. R. Jelenkovic and A. A. Lazar, "Asymptotic results for multiplexing subexponential on-off processes," *Advances in Applied Probability*, 1998 (forthcoming).

- 13. D. Bertsekas and R. Gallager, *Data Networks*, 1st ed., Prentice-Hall, Englewood Cliffs, N. J., 1987, pp. 114–122.
- N. Kahale and P. E. Wright, "Dynamic global packet routing in wireless networks," *Proc. IEEE INFOCOM '97*, Kobe, Japan, Apr. 1997, pp. 1416–1423.
- J. C. Chuang and N. Sollenberger, "Dynamic packet assignment for advanced Internet cellular services," *Proc. IEEE GLOBECOM '97, Phoenix, Ariz.*, Nov. 3–8, 1997, pp. 1596–1600.

(Manuscript approved September 1998)

MATTHEW ANDREWS received a B.A. in mathematics



from Oxford University in the United Kingdom and a Ph.D. in theoretical computer science from the Massachusetts Institute of Technology in Cambridge. A member of technical staff in the

Mathematics of Networks and Systems Department at Bell Labs in Murray Hill, New Jersey, Dr. Andrews works on wireless resource management, packet scheduling, and network design.

SIMON C. BORST holds an M.S. from the University of



Twente in Enschede, The Netherlands, and a Ph.D. from Tilburg University in The Netherlands, both in applied mathematics. A member of technical staff in the Mathematics of Networks and Systems

Research Department at Bell Labs in Murray Hill, New Jersey, Dr. Borst works on mathematical analysis of resource allocation problems in communication networks and computer systems.

FRANCIS DOMINIQUE is a member of technical staff in



the Base Station and Radio Department of the Wireless Networks Group in Whippany, New Jersey, where he is currently working on the analysis, design, and development of various digital signal processing subsystems

for third-generation CDMA wireless systems. Mr. Dominique received a B.S. in electronics and communications engineering from Pondicherry Engineering College in India and an M.S. in electrical engineering from Virginia Polytechnic Institute and State University in Blacksburg. PREDRAG R. JELENKOVIC, a member of technical staff in



the Mathematics of Networks and Systems Department at Bell Labs in Murray Hill, New Jersey, received a Ph.D. in electrical engineering from Columbia University, New York City. Dr. Jelenkovic works on the math-

ematical analysis of various queuing models, with particular interest in subexponential traffic characteristics.

KRISHNAN KUMARAN, who works on modeling, analy-



sis, and simulation of resource management issues in wireless networks and ATM, received a B.Tech. in mechanical engineering from the Indian Institute of Technology in Madras and a Ph.D. in physics from

Rutgers University in New Brunswick, New Jersey. Dr. Kumaran is a member of technical staff in the Mathematics of Networks and Systems Department at Bell Labs in Murray Hill, New Jersey.

K. G. RAMAKRISHNAN is a distinguished member of



technical staff in the Mathematics of Networks and Systems Department at Bell Labs in Murray Hill, New Jersey, where he works on mathematical optimization, integer programming, routing in ATM net-

works, and wireless applications. He earned a B.S. in electrical engineering from the Indian Institute of Technology in Kanpur and M.S. and Ph.D. degrees in computer science from the State University of Washington in Seattle.

PHILIP A. WHITING earned a B.A. in mathematics



from Oxford University in the United Kingdom, an M.Sc. in probability and statistics from the University of London in the United Kingdom, and a Ph.D. in engineering from the University of Strathclyde

in Glasgow, Scotland. As a member of technical staff in the Mathematics of Networks and Systems Department at Bell Labs in Murray Hill, New Jersey, Dr. Whiting conducts research on mathematical models of wireless systems, including teletraffic models and information theory.