# Power Management for Energy-Aware Communication Systems

CURT SCHURGERS, VIJAY RAGHUNATHAN, and MANI B. SRIVASTAVA
University of California, Los Angeles

System-level power management has become a key technique to render modern wireless communication devices economically viable. Despite their relatively large impact on the system energy consumption, power management for radios has been limited to shutdown-based schemes, while processors have benefited from superior techniques based on dynamic voltage scaling (DVS). However, similar scaling approaches that trade-off energy versus performance are also available for radios. To utilize these in radio power management, existing packet scheduling policies have to be thoroughly rethought to make them energy-aware, essentially opening a whole new set of challenges the same way the introduction of DVS did to CPU task scheduling. We use one specific scaling technique, dynamic modulation scaling (DMS), as a vehicle to outline these challenges, and to introduce the intricacies caused by the nonpreemptive nature of packet scheduling and the time-varying wireless channel.

Categories and Subject Descriptors: H.1.1 [**Models and Principles**]: Systems and Information Theory—*General systems theory*; C.2.1 [**Computer-Communication Networks**]: Network Architecture and Design—*Wireless communication*

General Terms: Algorithms, Performance, Design

Additional Key Words and Phrases: Energy-efficient, wireless communications, adaptive, scaling

## 1. INTRODUCTION

### 1.1 System Power Management

In commercial and research efforts alike, recent trends have shifted the emphasis in system design away from ever-increasing throughput alone. Instead, energy and power consumption are becoming ever more formidable and important constraints to address [Benini and de Micheli 1997; Pedram 2001]. Indeed,

technology integration and miniaturization have intensified cooling and packaging challenges. In addition, the energy supply itself is often limited to build-in batteries, as untethered operation is desired for mobility and portability. A class of systems where this observation holds especially true is that of personal communication devices, which are also becoming increasingly prevalent.

To reduce power spending without sacrificing performance, system designers have drawn on the concept of "energy or power awareness." The key idea is that the system only delivers the performance that is strictly required, thereby avoiding superfluous power consumption. Power and energy awareness are therefore often paraphrased as "the right power/energy at the right time and the right place." When designing according to this concept, there are two aspects to consider: the technique itself that introduces the awareness, and the power management strategy that exploits it.

The oldest and most straightforward technique is to shut down unused parts of the system [Benini et al. 1999]. Shutdown-based power management has been explored for hard disks, displays [Lorch and Smith 1998], and communication modules [Wang and Mandayam 2001], among others. For processors, it has been incorporated in the kernel of real-time operating systems (RTOS) [Srivastava et al. 1996]. However, a breakthrough came with the development of dynamic voltage scaling (DVS), a technique for digital circuits that is more effective than shutdown [Chandrakasan et al. 1992]. The reason is the convex nature of the power–speed curve, which comes about by varying the operating voltage. Numerous DVS-based power management strategies have been proposed to harness this potential [Burd et al. 2000; Govil et al. 1995; Gruian 2001; Gutnik and Chandrakasan 1997; Krishna and Lee 2000; Manzak and Chakrabarty 2000; Nielsen et al. 1994; Shin and Choi 1999; Weiser et al. 1994; Yao et al. 1995].

Unfortunately, energy awareness brought forth by DVS is restricted to digital circuits, such as processors, which can leverage DVS in their task-scheduling engine. Voltage scaling not being applicable; is it nevertheless possible to find a similar technique for other parts of the system? Or is shutdown the best we can do?

## 1.2 Power Management in Communication Subsystems

The search for superior power management techniques is especially relevant for the communication subsystem, and in this paper, we introduce and discuss a powerful scaling-based approach that can be viewed as the counterpart of DVS.

The importance of radio power management arises from the fact that communication is often the dominant power hog in the complete system [Raghunathan et al. 2002]. While the inherent energy consumption of digital circuits is rapidly decreasing due to Moore's law and ingenious design techniques, the power that is radiated to carry information does not follow this trend. This observation is particularly true for wireless RF (radio-frequency) devices, where the radiated power depends on the transmit distance and is physically constrained by Maxwell's laws. With the increasing proliferation

of cellular phones with Internet browsing capabilities, Bluetooth-connected PDAs, and wireless LANs, radio power management will undoubtedly gain in importance.

Radios typically posses a number of control knobs, which gives rise to a convex power-performance curve. As such, these knobs can be used to introduce energy awareness in a similar fashion as the operating voltage in DVS. Possible radio control knobs are the modulation level, the error coding, or combinations, or both. We refer to the resulting techniques as dynamic modulation scaling (DMS), dynamic code scaling (DCS), and dynamic modulation-code scaling (DMCS). To provide energy awareness, they need to be integrated into the system power management.

The aim of this paper is to provide insight into such scaling-based radio power management strategies, the associated challenges, and possible solutions. Just as people created energy-aware versions of RTOS task-scheduling policies, we investigate energy-aware packet scheduling as part of these power management strategies. As a vehicle for this exploration, we focus on dynamic modulation scaling (DMS), which we first introduced in Schurgers et al. [2001a]. It is a readily accessible control knob in a number of existing communication systems.

## 2. RELATED WORK

DMS for radios can be viewed as the counterpart of DVS for digital circuits, as both exploit the presence of a convex power–speed curve via the notion of scaling. We can therefore find inspiration in the way task scheduling was extended to make it energy-aware, when developing energy-aware packet-scheduling policies. Nevertheless, a direct correspondence is impossible due to the inherent differences between radios and digital circuits, which we detail in Section 3.4. Voltage scaling was first included in self-timed circuits [Nielsen et al. 1994] and later extended toward synchronous ones [Gutnik and Chandrakasan 1997], where a buffer is used to smooth the load and steer the adaptation. The initial DVS work in operation system research was in the context of a workstation like environment [Govil et al. 1995; Weiser et al. 1994], with average throughput as the performance metric. Task scheduling under real-time constraints, that is, the presence of deadlines, is considered in Burd et al. [2000], Gruian [2001], Krishna and Lee [2000], Manzak and Chakrabarty [2000], Shin and Choi [1999], and Yao et al. [1995].

In the realm of communications, a shutdown approach is proposed in Wang and Mandayam [2001], leveraging the time-varying nature of the wireless channel. Transmissions are deferred to times where the channel can support energy-efficient data transmission, while taking into account various delay constraints. A scheduling technique based on code scaling (DCS) is described in Prabhakar et al. [2001]. It finds an energy-efficient schedule for a set of packets that have to be transmitted before an overall deadline, and is similar to the work on processor scheduling presented in Yao et al. [1995]. The basic concept of modulation scaling (DMS) was first proposed by Schurgers et al. [2001a]. Here, we build upon this concept to describe a complete framework for dynamic power

management of the radio communication subsystem in wireless embedded systems.

We explore radio power management through DMS, which essentially relies on the ability of the radio to change its modulation on the fly. This is indeed practically feasible, and appropriate hardware architectures have been developed [Cho and Samueli 2000]. Originally, these architectures were used for a different purpose, namely adapting the modulation to maximize the system throughput [Balachandran et al. 1999; Ue et al. 1998; Webb and Steele 1995]. In this case, the appropriate choice of the modulation level only depends on the current condition of the wireless channel, and does not involve any scheduling decisions. On the other hand, as we illustrate in this paper, energy awareness is strongly related to scheduling, where current decisions and future events are tightly interwoven. Although utilizing the same radio control knob, DMS and throughput maximization not only serve a different purpose, but also require completely different adaptation policies. This change in mindset is similar to the one that occurred in the CPU world when the design goal was switched from throughput maximization to energy efficiency.

## 3. DYNAMIC MODULATION SCALING (DMS)

When considering energy-aware techniques beyond shutdown, such as DMS, we require some more insight into the operation of radios. In this section, we introduce the basics of modulation scaling.

### 3.1 Energy and Delay Breakdown

The first step in understanding DMS is investigating how energy and throughput depend on the modulation level in a digital wireless communication system. In order to transmit information, bits are coded into channel symbols, which correspond to different waveforms [Proakis 1995]. The number of possible waveforms determines how many bits are coded into one symbol, which is given by the modulation level $b$, expressed in number of bits per symbol. The average time to transmit 1 bit over the channel is the inverse of the average bit rate $R_b$, and is given by (1), where $R_S$ is the symbol rate (number of symbols that are transmitted per second).

$$T_{\text{bit}} = \frac{1}{R_b} = \frac{1}{b \cdot R_s}. \tag{1}$$

The energy associated with the transmission of 1 bit can be expressed as (2). The power needed to generate the information carrying electro-magnetic waves is delivered mainly by the power amplifier, and is denoted by the *transmit power* $P_S$. The remainder of the power consumption of the radio is lumped into $P_E$, the *electronics power*.

$$E_{\text{bit}} = [P_S + P_E] \cdot T_{\text{bit}}. \tag{2}$$

We note that (2) represents the energy consumption of the transmitter device, but it can also be used for the receiver, by setting $P_S$ equal to zero. The value of $P_E$ is not necessarily the same as that of the transmitter, of course.

Table I. Scaling Function $f(b)$ and $\Gamma$ for Different Modulation Schemes

| | $2^b$-QAM | $2^b$-PSK | $2^b$-PAM |
|---|---|---|---|
| $f(b)$ | $2^b - 1$ | $\sin\left(\frac{\pi}{2^b}\right)^{-2}$ | $\frac{2^{2b}-1}{3}$ |
| $\Gamma$ | $1/3 \cdot \left[Q^{-1}\left(\left(1 - \frac{1}{2^{b/2}}\right)^{-1} \cdot \frac{b \cdot \text{BER}}{4}\right)\right]^2$ | $1/2 \cdot \left[Q^{-1}\left(\frac{b \cdot \text{BER}}{2}\right)\right]^2$ | $1/2 \cdot \left[Q^{-1}\left(\left(1 - \frac{1}{2^b}\right)^{-1} \cdot \frac{b \cdot \text{BER}}{2}\right)\right]$ |

Expressions (1) and (2) give the delay and energy associated with transmitting a bit over the wireless link, and therefore operate on the level of the OSI physical layer [Tenenbaum 1990]. These bits do not simply represent the bare application information (so-called "useful bits"), but contain contributions from higher-layer overhead, such as headers, channel coding, training sequences, control packets, and so on [Lettieri et al. 1999]. As we focus here on DMS, we keep all other parameters and protocol settings unchanged. The energy and delay per "useful bit" are thus a linear function of their corresponding values on the physical layer, such that we can focus on (1) and (2). We note that energy-efficient techniques on the other layers are still perfectly applicable and beneficial in addition to DMS.

To analyze DMS, we need to derive the detailed relationship between the energy and the modulation level. The scheme that is probably most amendable to scaling, due to its ease of implementation and analysis, is quadrature amplitude modulation (QAM) [Balachandran et al. 1999; Ue et al. 1998; Webb and Steele 1995]. The resulting bit error rate (BER) is well approximated by (3) [Proakis 1995]. In this equation, $E_N$ is the noise energy per symbol, factor $A$ contains all transmission loss components, and $\eta$ is the linearized efficiency of the power amplifier. The $Q(\cdot)$ function is defined as (4).

$$\text{BER} = \frac{4}{b} \cdot \left(1 - \frac{1}{2^{b/2}}\right) \cdot Q\left(\sqrt{3 \cdot \frac{A \cdot \eta}{2^b - 1} \cdot \frac{P_S}{E_N \cdot R_S}}\right) \qquad (3)$$

$$Q(x) = \frac{1}{\sqrt{2\pi}} \cdot \int_x^\infty \exp\left(-\frac{1}{2} \cdot t^2\right) dt = \frac{1}{2} \cdot \text{erfc}\left(\frac{x}{\sqrt{2}}\right). \qquad (4)$$

By solving for the transmit power, we obtain (5), where parameter $C_S$ is defined as (6). The expressions for $f(b)$ and $\Gamma$ are listed in the first column of Table I.

$$P_S = C_S \cdot f(b) \cdot R_S \qquad (5)$$

$$C_S = \frac{E_N}{A \cdot \eta} \cdot \Gamma. \qquad (6)$$

$E_N$ is a function of the receiver implementation and the operating temperature. $A$ depends on distance and the propagation environment, and can vary with time. Neither of them varies with $b$. In addition, due to the $Q^{-1}(.)$ function, $\Gamma$ is only very weakly dependent on $b$. Although we do take this dependency into account in the simulations, $C_S$ is thus approximately a constant when we scale the modulation (i.e., vary $b$). The main benefits from modulation scaling are due to $f(b)$.

Fig. 1. Fractional modulation-level settings.

The power consumption of the electronic circuitry, which is largely analog, $P_E$, can be written as (7) [Cho and Samueli 2000; Schurgers et al. 2001a]. $C_E$ is a constant that depends on the radio architecture, the circuit implementation, and the semiconductor technology.

$$P_E = C_E \cdot R_S. \tag{7}$$

With (5) and (7), expression (2) becomes an explicit function of the modulation level:

$$E_{\mathrm{bit}} = [C_S \cdot f(b) + C_E] \cdot \frac{1}{b}. \tag{8}$$

The ratio of parameters $C_S$ and $C_E$ can thus be viewed as an indication of the relative importance of the transmit power versus the electronics power. In a similar way, we can derive expressions for phase shift keying (PSK) and pulse amplitude modulation (PAM) [Proakis 1995]. With the appropriate definitions of $f(b)$ and $\Gamma$ as in Table I, (8) remains valid. In general, DMS is applicable to other scalable modulation schemes as well.

## 3.2 Granularity Effects

In the previous section, we implicitly assumed the modulation level could be varied continuously. However, strictly speaking, the expressions in Table I are only valid for integer values of $b$. In the case of QAM, they are exact only for even integers, but are reasonable approximations when $b$ is odd (so-called "nonsquare" constellations) [Proakis 1995].

Furthermore, it is impractical to change the modulation level at arbitrary time instants, since both sender and receiver need to know the exact modulation scheme that is used. This requires a kind of negotiation between the two, resulting in protocol overhead. It makes sense to limit the rate of adaptation, for example, by restricting it to periodic instants or the start of packet transmissions. The implication toward implementation is that the receiver has to be designed such that it can appropriately reconfigure its processing for the right modulation level. This is one crucial difference between modulation scaling and voltage scaling, which, as we will see later, also affects the power management strategies.

However, in between two modulation updates, we can define a fractional modulation level. For example, the first half of the packet could be sent with modulation $b_1$ and the second half with $b_2$, see Figure 1. As a result, the average energy and delay per bit are a linear interpolation between the corresponding values of the two modulation levels, as is apparent from (9) and (10).

$$E_{\mathrm{bit}} = \frac{E_{\mathrm{packet}}}{L} = E_{\mathrm{bit}}^1 \cdot \beta + E_{\mathrm{bit}}^2 \cdot (1 - \beta) \tag{9}$$

Fig. 2.   Energy-delay trade-off for QAM.

Table II.  Simulation Settings

| | |
|---|---|
| $R_S$ | 1 MH$_Z$ |
| BER | $10^{-5}$ |
| $C_S$ (4-QAM) | 12 nJ |
| $C_E$ | 15 nJ |

$$T_{\text{bit}} = \frac{T_{\text{packet}}}{L} = T_{\text{bit}}^1 \cdot \beta + T_{\text{bit}}^2 \cdot (1 - \beta). \tag{10}$$

Each modulation scheme has a $b_{\min}$, which is the minimum practically achievable modulation. If we are allowed even more delay per bit, we can shutdown the radio for a while ($b = 0$). However, as mentioned before, shutdown does not reduce the energy per bit. The maximum modulation $b_{\max}$ is only bounded by implementation choices and maximum transmit power. Between $b_{\min}$ and $b_{\max}$, the modulation can be scaled with granularity $\delta$, which is a design choice. This minimum and maximum level and granularity effect are similar to what is encountered in voltage scaling as well.

Figure 2 illustrates the energy-delay trade-off for QAM, with the numerical values from Table II. $C_S$ varies slightly with the modulation level, and the value at operating point $b = 2$ is listed in the table. The curve labeled "ideal" corresponds to (1) and (8). The circles indicate constellations that can be realized, and the solid line gives the values that are obtained through the interpolation we just explained. For QAM, $b_{\min}$ is equal to 2. We see that we can use (1) and (8) for analysis purposes as very tight approximations to what is practically realizable.

$E_{bit}$ (μJ)



Fig. 3.   Energy-delay trade-off for different $C_S$.

## 3.3 Region of Modulation Scaling

Figure 3 shows the same trade-off for QAM for different values of $C_S$. The values of $C_S$ shown in the figure are for operating point $b = 2$, that is, 4-QAM, and correspond to $P_S$ varying from 2.25 mW to 144 mW. The other settings are those of Table II. As $C_E$ is kept constant, a higher value of $C_S$ thus indicates that the transmit power portion becomes more dominant compared to the electronics power portion.

DMS essentially utilizes the effect that, by varying the modulation level, energy can be traded off versus delay. At the left side of the figure, lowering $b$ reduces the energy, at the cost of an increased delay. Alternatively, the power versus speed is convex in this region. Scaling beyond the point of minimum energy clearly does not make sense, as both energy and delay would increase. The operating region of DMS therefore corresponds to the portion of the curves to the left of their minimum energy point. It is clear that in this region, scaling is superior to shutdown, as the metric $E_{bit}$ is the total energy per bit, which will not change if the transmission is followed by a period of shutdown. We can also verify that the energy–delay curves are convex, such that a uniform stretching of the transmissions is the most energy-efficient (due to Jensen's inequality).

From Figure 3, we see that DMS is more useful for situations where $C_S$ is large, or in other words, where the transmit power dominates the electronics power. This is true except for communication systems with a very short transmission range. Also note that the electronics power behaves similarly to the leakage current in digital circuits. As leakage becomes more dominant with shrinking CMOS device dimensions [Sinha and Chandrakasan 2000], DVS will be faced with similar issues of operating region, as observed here.

Table III. Modulation Bounds

| $b_{min}$ (bits/symbol) | $b_{max}$ (bits/symbol) | $\delta$ (bits/symbol) |
|---|---|---|
| 2 | 8 | 0.5 |

## 3.4 Difference Between DMS and DVS

Despite the analogies between DVS and DMS, there are two important differences. First, a change in modulation level requires communication between sender and receiver. As explained in Section 3.2, the sender cannot decide on a new modulation setting midway through the packet transmission. Packet scheduling therefore exhibits this inherent time-granularity effect and can only be nonpreemptive. However, the exact packet size is known at the start of the transmission, in contrast to the task execution time in processors [Raghunathan et al. 2001]. We detail these points in Section 4.1.

Second, the wireless channel may vary over time, which means that the factor $A$ in (6) fluctuates. These variations have to be taken into account in the packet-scheduling engine. This is akin to a processor where capacity (in MIPS) would change over time, for example, due to interrupt handling. This phenomenon is indeed complex, but is beyond the designers' control in radio power management.

## 4. ENERGY-AWARE PACKET SCHEDULING

Just like DVS has driven energy-aware task scheduling beyond shutdown-based approaches [Burd et al. 2000; Govil et al. 1995; Gruian 2001; Gutnik and Chandrakasan 1997; Krishna and Lee 2000; Manzak and Chakrabarty 2000; Nielsen et al. 1994; Shin and Choi 1999; Weiser et al. 1994; Yao et al. 1995], DMS paves the way toward new energy-aware packet-scheduling policies. The body of literature dealing with packet scheduling is vast, and, in principle, is suitable to be extended toward energy-aware versions using DMS [Demers et al. 1989; Lu et al. 1997; Parekh and Gallager 1994]. However, many challenges lie ahead, since such radio power management must deal with both traffic load and channel variations. In the next subsections, we describe two basic scheduling approaches to illustrate some of the challenges. They each highlight one of two different issues, namely the presence of deadlines and channel variations.

In all simulations reported in this section, a special field in the packet header, encoded with 4-QAM, is used to communicate the modulation level for the rest of the packet to the receiver. We chose this option for its simplicity and fault tolerance, assuming we have control over the protocol stack. If this is not available, an alternative is to use separate control packets. The DMS scheduler, residing at the link-layer, might still require channel quality information from the lower layers, similar to other adaptive protocols [Balachandran et al. 1999; Lettieri et al. 1999; Thoen et al. 2000]. Table III defines the possible values, which can be encoded into 4 bits. This overhead for each packet is incorporated in all simulation results.

## 4.1 Real-Time Scheduling in a Time-Invariant Channel

In this section, we consider the problem of scheduling multiple real-time traffic streams, while the wireless channel does not vary in time. In each stream, packets arrive periodically and have a deadline by which they have to be sent. This is a valid model for multimedia streams. We choose the deadline of each packet to be the arrival time of the next packet in that stream. The size of the individual packets might vary (e.g., in MPEG video or perceptual audio codecs), but there is a maximum known packet size for each stream. As mentioned before, the packet scheduler should not interrupt an ongoing transmission such that the policy has to be nonpreemptive. This is a new constraint compared to the energy-aware task scheduling policies for RTOS, which are all preemptive (i.e., tasks can be suspended and resumed later on) [Burd et al. 2000; Gruian 2001; Krishna and Lee 2000; Manzak and Chakrabarty 2000; Shin and Choi 1999; Yao et al. 1995]. The setup we consider here is analogous to the DVS work described in Raghunathan et al. [2001]. Besides the preemptive nature, the scheduling algorithm of Raghunathan et al. [2001] deals with task-length variability in a stochastic fashion, as the exact execution time is not known at the start of the task. However, in packet scheduling, the exact packet length is known at the start of the transmission.

A condition that guarantees schedulability for nonpreemptive scheduling is derived in Jeffay et al. [1991]. When it is satisfied, earliest-deadline-first scheduling (EDF) always results in a valid schedule. An optimal energy-aware scheduling routine is too computationally intensive as the problem is NP-complete; but we propose a practical algorithm, which consists of two steps [Schurgers et al. 2001b].

1. *Admission step*: When a new stream is admitted to the system, we calculate a static scaling factor, $\alpha_{\text{static}}$, assuming all packets are of maximum size. This factor is the minimum possible such that if the modulation setting for each packet were scaled by it, the schedulability test would still be satisfied. In other words, it computes the slowest transmission speed at which all the packet streams are just schedulable.

2. *Adjustment step*: At run-time, packets are scheduled using EDF. Before transmission starts, the actual size of each packet is known. We calculate an additional scaling factor, $\alpha_{\text{dyn}}$, such that the transmission finishes when that of a maximum size packet would have. Since step 1 assumes the maximum packet size, the schedulability is guaranteed. If the system were still idle after the packet transmission, we would stretch the transmission until the packet's deadline or the arrival time of a new packet (which is known due to the periodic nature of the traffic). This extra scaling factor is called $\alpha_{\text{stretch}}$.

The scheduler combines all three scaling factors to get the overall modulation that is used for the current packet [Schurgers et al. 2001b]. To evaluate the performance of our scheme, we have carried out a number of simulations. The basic settings and modulation settings are those of Tables II and III, respectively. The

Table IV.  Number of Streams with Given Period

|          | 20 ms | 25 ms | 40 ms | 50 ms |
|----------|-------|-------|-------|-------|
| $U = 0.82$ | 8 | 6 | 4 | 4 |
| $U = 0.77$ | 8 | 6 | 2 | 4 |
| $U = 0.74$ | 8 | 4 | 4 | 4 |
| $U = 0.69$ | 6 | 4 | 6 | 4 |



Fig. 4.   Relative energy for real-time energy-aware packet scheduling.

size of each packet is independently chosen in a uniformly random fashion between the maximum packet size $L_{max}$ of 8000 bits and a minimum value $L_{min}$, which is a parameter we vary in our simulations. We consider four different scheduling scenarios. For each of them, a row in Table IV lists the number of streams with a given period and the resulting total link utilization when all packets are of maximum size. For example, in the fourth scenario in Table IV, there are six streams with a period of 20 ms, four with a period of 25 ms, six with a period of 40 ms, and four with a period of 50 ms, resulting in a link utilization $U$ of 0.69 when $L_{min} = L_{max}$.

Figure 4 plots the energy consumption of our scheduling scheme, normalized against one without scaling ($b = b_{max}$ at all times). For $U = 0.82$, we have separated the effect of the different scaling factors. When only using $\alpha_{static}$, the transmissions are slowed down uniformly without exploiting the run-time packet length variations. These are leveraged through $\alpha_{dyn}$, where the energy decreases as the packet size variation increases ($L_{min}$ decreases). The effect of $\alpha_{strech}$ is marginal here. For the other utilizations, we only show the results when combining all three scaling factors. As expected, more energy savings are achieved when the utilization is lower.

The power management scheme described here, essentially exploits traffic load variations on two levels to introduce energy awareness.

1. Variations in overall utilization are handled by the admission step in our algorithm through $\alpha_{\text{static}}$. These are due to changes in number of streams, which are likely to occur over relatively large time scales.
2. Variations in individual packet sizes on the other hand occur at much smaller time scales. These cannot be handled during admission, but are exploited in the adjustment step of our algorithm.

## 4.2 Nonreal-Time Scheduling in Time-Variant Channel

In this section, we highlight a different issue in radio power management, namely the effects of time variations in the wireless channel. The closest equivalent in CPU task scheduling is a time-varying processor capacity. To introduce some of the challenges, we consider a rudimentary scenario: the transmission of a single data stream that has no hard deadline associated with it, but only an average data rate constraint. This model is useful in the case of a file transfer, for example. In Schurgers and Srivastava [2002], we also illustrate how this specification can be used to provide a soft real-time constraint.

As discussed in Section 3.1, the transmission loss A captures the effect of the wireless channel. In the presence of time variations, this factor is split up into two components as in (11), where $\overline{A}$ represents the long-term average value and $\alpha$ contains the normalized time variations. The behavior of the gain factor $\alpha$ can be characterized by two statistics: a probability density function and a Doppler rate, which describes the time correlation [Jakes 1994; Proakis 1995].

$$A = \overline{A} \cdot \alpha \tag{11}$$

In all techniques that adapt a radio parameter to channel variations, an estimate of the current channel condition is needed [Balachandran et al. 1999; Lettieri et al. 1999]. This is obtained through channel estimation, which is updated regularly. The update rate $f_{\text{update}}$ is chosen such that the channel remains approximately constant between updates, yet the overhead of the estimation process is limited. In addition, predictive compensation techniques can be applied [Thoen et al. 2000]. Although deciding the appropriate update rate, as a compromise between overhead and performance degradation, is an important issue, a detailed study falls outside the scope of this paper.

In the previous section, DMS turned into a scheduling problem because of the interaction between multiple streams. Here, we only have one stream, but the presence of a time-varying channel again makes the choice of the best value of b a scheduling issue. The current decision critically depends on how good or bad the channel will be, that is, whether it is more energy efficient to send now or later. However, in the specific scenario we consider here, namely that the location of the average throughput is the only additional constraint, the problem can be greatly simplified. In Schurgers and Srivastava [2002], we prove that there exists a set of thresholds $d_i$ that directly link the current channel condition to the optimal choice of $b$. Equation (12) presents a direct generalization of the results presented in Schurgers and Srivastava [2002]. It is not valid for very

high values of $C_E/C_S$, but holds the ones chosen here [Schurgers 2002]. We refer to the resulting energy-efficient DMS packet-scheduling policy as "loading in time," as it was inspired by adaptive bit-loading techniques for a class of radio systems called multicarrier systems [Chow et al. 1995; Hughes-Hartogs 1987]. Note that only those values of $b$ are assigned that did not result from the interpolation of Section 3.2. The reason why they are not included in the optimal assignment is that the energy versus delay curve is no longer strictly convex there [Schurgers 2002].

$$\begin{cases} 0 \le \alpha < d_1 \Rightarrow b = 0 \\ d_i \le \alpha < d_{i+1} \Rightarrow b = b_{\min} + (i-1), \quad i = 1..(K-1) \\ d_K \le \alpha < \infty \Rightarrow b = b_{\max} \end{cases}$$

$$K = 1 + b_{\max} - b_{\min}. \tag{12}$$

Furthermore, the thresholds $d_i$ are mutually related according to (13) and (14) [Schurgers 2002]. The fact that both transmit and electronics power are zero when $b = 0$ is taken into account here.

$$d_i = \max[\theta \cdot 2^{i-1}, d_1] \qquad i = 2..K \tag{13}$$

$$\theta = b_{\min} \cdot \left[ \frac{2^{b_{\min}} - 1}{d_1} + \frac{C_E}{C_S} \right]^{-1} \cdot 2^{b_{\min}-1}. \tag{14}$$

There is only one independent parameter left, which can be solved from the constraint on the desired average data rate $b_{av}$, expressed in average number of bits per symbol. This constraint can be written as (15), where $G(x)$ is defined in (16) [Schurgers and Srivastava 2002]. Here, $F(x)$ is the cumulative distribution function of the gain factors $\alpha$.

$$b_{av} = b_i \cdot G(d_1) + \sum_{i=2}^{K} G(d_i) \tag{15}$$

$$G(x) = 1 - F(x). \tag{16}$$

The thresholds thus only depend on the statistics of the wireless channel, which can be estimated online. We no longer have to know the exact behavior of the channel over time to achieve the energy-optimal scheduling policy. Figure 5 shows the simulated performance of this radio power management scheme versus the average throughput constraint. As before, we used the values of Tables II and III. The channel exhibits the correlated Rayleigh fading, such that $G(x)$ is given by (17). This corresponds to the case where there is no line-of-sight path between sender and receiver, which is the predominant model used in literature [Proakis 1995]. The time correlation of the channel is characterized by a Doppler rate of 50 Hz, and simulated as proposed in Jakes [1994].

$$G(x) = \exp(-x) \tag{17}$$

We selected an update rate $f_{update}$ of 1 kHz for channel estimation, the overhead of which is included in the reported simulation results. The maximum

$E_{bit}$ (µJ)



Fig. 5. Energy versus average bit rate for time-varying channel (Rayleigh fading).

possible transmit power is 1 W. Curve 1 in Figure 5 plots the behavior of the "loading in time" scheduling policy that we described in this section. It is superior to scaling with "constant b" (curve 2), where the modulation is uniformly slowed down based on the average throughput, but channel variations are not taken into account. The difference between curves 2 and 3, which shows the same uniform scaling in a nontime-varying channel, illustrates the performance degradation associated with channel variations. Beyond $b_{\min} = 2$ bits/symbol, we resort to shutdown, and both these curves flatten out, which is as expected from our discussion in Section 3.2. However, curve 1 keeps on decreasing when lowering $b_{av}$, and can even outperform scaling in a non-time-varying channel (curve 3). The reason is that we still use shutdown, but only the very best time intervals (with $\alpha > 1$) are selected to send information. For curve 2, the shutdown was periodic, without taking the channel into account.

Finally, we also compare the performance to a scheme that is not energy-aware, but tries to achieve a "maximum throughput" possible. In this case, b is adapted to yield its maximum value without violating the maximum transmit power [Balachandran et al. 1999; Ue et al. 1998; Webb and Steele 1995]. As this is only based on the current channel condition, scheduling issues never arise. The benefits of energy-awareness, where a reduced throughput requirement is leveraged to yield energy savings, are again substantial.

## 5. CONCLUSIONS

Energy efficiency is gaining importance as a system design consideration, especially in portable communication devices. Radio-level power management based on shutdown simply turns off the radio when it is not used. However,

scaling-based techniques have the potential to be vastly superior in terms of energy savings, but necessitate more intricate power management. Packet scheduling has to be rethought to include energy-awareness, similar to the way voltage scaling prompted the search for energy-aware task scheduling. Numerous wireless packet scheduling techniques exist, which are all candidates to make it into an energy-aware version.

In this paper, we focused on dynamic modulation scaling, DMS, as just one of the possible radio control knobs that can be used for scaling. Even in the elementary scenarios we considered here, it illustrates the intricacies and challenges of energy-aware packet scheduling. These will be compounded when stringent delay/throughput constraints have to be satisfied in the presence of a time-varying wireless channel, which is inherently stochastic in nature. It is expected that resulting policies will not be able to guarantee service, but be characterized by reliability bounds, where the tightness of the bounds can be traded off for energy savings. Radio power management therefore has to take both traffic and channel aspects into account simultaneously. Other research challenges that remain are enhancing these strategies to also handle multiple devices on a shared channel. This likely requires a distributed manipulation of the radio control knobs, by incorporating it in the medium access control layer.

The resulting radio power management eventually has to be integrated into a system-wide solution, where energy and latency can be traded off across the computation–communication subsystem boundaries. The overall vision is thus a coordinated power management, intermixing both task and packet scheduling policies in a networked system. Furthermore, these ideas, while presented here in the context of radios, are potentially generalizable to wired communications as well, which may offer other control knobs for power-speed scaling.

REFERENCES

BALACHANDRAN, K., KADABA, S. R., AND NANDA, S. 1999. Channel quality estimation and rate adaptation for cellular mobile radio. *IEEE Journal on Selected Areas in Communications 17*, 7, 1244–1256.

BENINI, L. AND DE MICHELI, G. 1997. *Dynamic Power Management: Design Techniques and CAD Tools*. Kluwer, Norwell, MA.

BENINI, L., BOGLIOLO, A., AND DE MICHELI, G. 1999. A survey of design techniques for system-level dynamic power management. *IEEE Transactions on VLSI Systems 8*, 3, 813–833.

BURD, T., PERING, A., STRATAKOS, A., AND BRODERSEN, R. 2000. A dynamic voltage scaled microprocessor system. *IEEE Journal of Solid-State Circuits 35*, 11, 1571–1580.

CHANDRAKASAN, A., SHENG, S., AND BRODERSEN, R. 1992. Low-power CMOS digital Design. *IEEE Journal of Solid-State Circuits 27*, 4, 473–484.

CHO, K. AND SAMUELI, H. 2000. A 8.75-MBaud single-chip digital QAM modulator with frequency-agility and beamforming diversity. In *Proceedings CICC'00*. Orlando, FL (May), 27–30.

CHOW, P., CIOFFI, J., AND BINGHAM, J. 1995. A practical discrete multitone transceiver loading algorithm for data transmission over spectrally shaped channels. *IEEE Trans. on Communications 43*, 2, 773–775.

DEMERS, A., KESHAV, S., AND SHENKER, S. 1989. Analysis and simulation of a fair queueing algorithm. *Computer Communication Review 19*, 4, 1–12.

GOVIL, K., CHAN, E., AND WASSERMAN, H. 1995. Comparing algorithms for dynamic speed-Setting of a low-power CPU. In *Proceedings MobiCom'95*. Berkeley, CA (Nov.), 13–25.

GRUIAN, F. 2001. Hard real-time scheduling for low energy using stochastic data and DVS processor. In *Proceedings ISLPED'01*. Huntington Beach, CA (Aug.), 46–51.

GUTNIK, V. AND CHANDRAKASAN, A. 1997. Embedded power supply for low-power DSP. *IEEE Transactions on VLSI Systems 5*, 4, 425–435.

HUGHES-HARTOGS, D. 1987. Ensemble modem structure for imperfect transmission media. U.S. Patents, nos. 4,679,227 (July 1987), 4,731,816 (March 1988), 4,833,796 (May 1989).

JAKES, W. 1994. *Microwave Mobile Communication*. John Wiley, New York.

JEFFAY, K., STANAT, D., AND MARTEL, C. 1991. On non-preemptive scheduling of periodic and sporadic tasks. In *Proceedings IEEE RTSS'91*. San Antonio, TX (Dec.), 129–139.

KRISHNA, C. AND LEE, Y. 2000. Voltage-clock-scaling adaptive scheduling techniques for low power in hard real-time systems. In *Proceedings RTAS'00*. Washington, DC (June), 156–165.

LETTIERI, P., SCHURGERS, C., AND SRIVASTAVA, M. 1999. Adaptive link layer strategies for energy efficient wireless networking. *Wireless Networks 5*, 5, 339–355.

LORCH, J. AND SMITH, A. 1998. Software strategies for portable computer energy management. *IEEE Personal Communications 5*, 3, 60–73.

LU, S., BHARGHAVAN, V., AND SRIKANT, R. 1997. Fair scheduling in wireless packet networks. *Computer Communication Review 27*, 4, 63–74.

MANZAK, A. AND CHAKRABARTY, C. 2000. Variable voltage task scheduling for minimizing energy or minimizing power. In *Proceedings ICASSP'00*. Istanbul, Turkey (June), 3239–3242.

NIELSEN, L., NIESSEN, C., SPARSØ, J., AND VAN BERKEL, K. 1994. Low power operation using self-timed circuits and adaptive scaling of the supply voltage. *IEEE Transactions on VLSI Systems 2*, 4, 391–397.

PAREKH, A. AND GALLAGER, R. 1994. A generalized processor sharing approach to flow control in integrated services networks: The multiple node case. *IEEE/ACM Transactions on Networking 2*, 2, 137–150.

PEDRAM, M. 2001. Power optimization and management in embedded systems. In *Proceedings ASP-DAC 2001*. Yokohama, Japan, 239–244.

PRABHAKAR, B., BIYIKOGLU, E., AND EL GAMAL, A. 2001. Energy-efficient transmission over a wireless link via lazy packet scheduling. In *Proceedings Infocom'01*. Anchorage, AK (April), 386–394.

PROAKIS, J. 1995. *Digital Communications*. McGraw-Hill Series in Electrical and Computer Engineering. McGraw-Hill New York.

RAGHUNATHAN, V., SCHURGERS, C., PARK, S., AND SRIVASTAVA, M. 2002. Energy-aware wireless sensor networks. *IEEE Signal Processing Magazine 19*, 2, 40–50.

RAGHUNATHAN, V., SPANOS, P., AND SRIVASTAVA, M. 2001. Adaptive power-fidelity in energy aware wireless systems. In *Proceedings RTSS'01*. London, UK (Dec.), 106–115.

SCHURGERS, C. 2002. Energy-aware wireless communications. Ph.D. dissertation, University of California at Los Angeles.

SCHURGERS, C., ABERTHORNE, O., AND SRIVASTAVA, M. 2001a. Modulation scaling for energy aware communication systems. In *Proceedings ISLPED'01*. Huntington Beach, CA (Aug.), 96–99.

SCHURGERS, C., RAGHUNATHAN, V., AND SRIVASTAVA, M. 2001b. Modulation scaling for real-time energy aware packet scheduling. In *Proceedings Globecom'01*. San Antonio, TX (Nov.), 3653–3657.

SCHURGERS, C. AND SRIVASTAVA, M. 2002. Energy efficient wireless scheduling: Adaptive loading in time. In *Proceedings WCNC'02*. Orlando, FL (Mar.), 706–711.

SHIN, Y. AND CHOI, K. 1999. Power conscious fixed priority scheduling for hard real-time systems. In *Proceedings DAC'99*. New Orleans, LA (June), 134–139.

SINHA, A. AND CHANDRAKASAN, A. 2000. Energy aware software. In *Proceedings VLSI Design 2000*. Calcutta, India (Jan.), 50–55.

SRIVASTAVA, M., CHANDRAKASAN, A., AND BRODERSEN, R.  1996.  Predictive system shutdown and other architectural techniques for energy efficient programmable computation. *IEEE Transactions. on VLSI Systems 4*, 1, 42–55.

TENENBAUM, A.  1990.  *Computer Networks*. Prentice-Hall, Englewood Cliffs, NJ.

THOEN, S., VAN DER PERRE, L., GYSELINCKX, B., ENGELS, M., AND DE MAN, H.  2000.  Predictive adaptive loading for HIPERLAN II. In *Proceedings VTC'00 Fall*. Boston, MA (Sept.), 2166–2172.

UE, T., SAMPEI, S., MORINAGA, N., AND HAMAGUCHI, K.  1998.  Symbol rate and modulation level-controlled adaptive modulation/TDMA/TDD System for High-Bit Rate Wireless Data Transmission. *IEEE Transactions on Vehicular Technology 47*, 4, 1134–1147.

WANG, H. AND MANDAYAM, N.  2001.  Delay and energy constrained dynamic power control. In *Proceedings Globecom'01*. San Antonio, TX (Nov.), 1287–1291.

WEBB, W. AND STEELE, R.  1995.  Variable rate QAM for mobile radio. *IEEE Transactions on Communications 43*, 7, 2223–2230.

WEISER, M., WELCH, B., DEMERS, A., AND SHENKER, B.  1994.  Scheduling for reduced CPU energy. In *Proceedings USENIX Symposium on Operating Systems Design and Implementation*. Monterey, CA (Nov.), 13–23.

YAO, F., DEMERS, A., AND SHENKER, S.  1995.  A scheduling model for reduced CPU energy. In *Proceedings 36th Annual Symposium on Foundations of Computer Science*. Milwaukee, WI (Oct.), 374–385.