Trade-off Between Power Consumption and Delay in Wireless Packetized Systems

Todd Coleman and Muriel Médard

Laboratory for Information and Decision Systems Massachusetts Institute of Technology Cambridge, MA 02139 colemant@mit.edu, medard@mit.edu

October 3, 2001

Abstract

In packetized wireless systems, coding allows correct reception of multiple packets colliding at a receiver. Thus data may not need to incur delays such as those due to backoff schemes in traditional ALOHA systems. However, there is a tradeoff between delay and power consumption. Recent work in this area has considered the case where multiple users are aware of the states of other users' queues. We consider a time-slotted multiple user system with random packet arrivals. The size of the packets and probability of arrival together represent the burstiness of the system. The time slots are considered to be long enough that capacity can be achieved over a single slot in a sense we define. We consider the difference in average power consumption when average delay is minimized, with and without knowledge of other users' queues. We also consider the case where average power is minimized without regard for delay. We present and analyze a simple scheme with limited information sharing about queues' states. Our scheme uses a broadcasttype code for the case of low queue lengths and a multiple-acces scheme in the case of large queue lengths. We show how this scheme allows trade-offs between power and delay.

1 Introduction

The performance of wireless nomadic data transfer systems can be characterized by a number of system qualities, including aggregate data rates, average bit transmission delay, and power consumption. Information theoretic considerations attempt to establish ultimate limits on reliable communication. Shannon capacity assumes that a steady stream of bits is to be transmitted at all times. Many data transfer systems, however, exhibit random packet arrivals. The size of the packets and probability of arrival together represent the burstiness of the system. This violates the assumption that bits are always available for transmission.

The time-slotted ALOHA system models systems with bursty arrivals. Its use is motivated by its simplicity: users attempt to transmit data as it arrives in their transmission buffers. If two or more users transmit at the same time, a collision occurs at the receiver. Traditional ALOHA systems require users to transmit packets without explicit coordination among users. In the event of a collision, packets are discarded and users retransmit the collided packets. The capacity of such systems has generally been analyzed in terms of packet throughput.

The stability of classical ALOHA systems has been studied extensively. For an infinite number of users, it has been found in [1] that the system is unstable. Stability regions have been found for systems with a finite number of users. To combat the instability, decentralized control schemes [2] and conflict resolution schemes [3] have been established. In general, such schemes attempt to avoid successive collisions by retransmitting with some backoff policy.

Modeling a collision as leading to loss of all packets at receivers is not always practical. The capture phenomenon, for instance, may yield correct reception of some portion of the data, for instance the data from coded slots. Moreover, users may be reliably received if, when transmitting, they take into account the worst case multiple access scenario that may arise. If coding can be implemented over sufficiently many bits, then users, when they transmit, may use the types of codes that achieve rates on the Cover-Wyner multiple access rate region.

Stability analysis of systems with multiple-packet reception capability in the presence of channel noise has been performed in [4]. The capacity region of such systems, in a sense we qualify later, that allows coding of packets *and* variable reliably received rates has recently been introduced [5]. It has been found that such a system's capacity region is the same as the capacity region of a multiple-access system where users continuously transmit. Furthermore, transmission policies that make use of detailed knowledge of users' queus, whether in a decentralized fashion or through centralized control such as a scheduler, do not improve capacity. Hence, capacity of such systems is in general independent of burstiness and queue information availability. Many coding schemes were shown to be optimal, ensuring long-term stability while achieving rates inside the Cover-Wyner region. The impact of burstiness and queue information, however, when considering delay and power consumption, was not illustrated in [5]. Clearly, queue information is not altogether useless. While it does not affect the type of capacity we consider, we would expect it to influence other performance parameters, such as delay.

An investigation of the trade-offs between minimum average power required to meet some quality of service cost (stringent delay or probability of buffer overflow) [6] has been performed. This analysis addressed bursty multi-user channels in the presence of fading. The investigation used centralized control to combat both fading and burstiness to deliver a stringent delay constraint.

We consider the difference in average power consumption when average bit delay is minimized to 0 under two different scenarios. We investigate a control scheme where the transmission policy of each user relies on detailed information about the amount of data in all users' queues. We also investigate a control scheme where each user's policy relies only on the amount of data in that particular user's queue. Hence, the impact of queue information sharing in the presence of burstiness is characterized. We also consider minimizing power consumption without regard for delay, which turns out to be infinite. Finally, we present and analyze a simple scheme that has some optimal long-term stability properties, combats burstiness and collisions by superimposing codes anticipating different levels of interference, and is sensitive to average bit delay. The system uses a multiple access channel [7] type code when all users' queue lengths are long. Otherwise, it uses a broadcast [8] type code. Instead of combating burstiness by performing probabilistic backoff policies after collisions, the system uses rate-splitting to achieve variable reliably received rates as a function of the uncertainty of other users'



Figure 1: The *M*-user ALOHA model.

presence. As user queue lengths become long, the system switches to coding in multipleaccess mode, where users code for each others' presence at optimal aggregate rates (as provided by the dominant face of the Cover-Wyner region for multiple access channels).

2 Model and Background

The multiple-access system we consider is illustrated in Figure 1. Our system differs from traditional time-slotted ALOHA systems in a number of ways:

- We model a multiple access where users share a single channel with no multiplicative attenuation but with additive white Gaussian noise (AWGN).
- Time slots are very long in terms of bits to be transmitted. Consequently, transmission data may be coded to achieve rates near information-theoretic bounds. We may also consider coding over several time slots, but note that this may complicate the discussion without providing much insight. This long time slot model may be particularly appropriate when dealing with channels with small signal-to-noise ratios, where turbo codes have been found to achieve rates near capacity with sustainable probability of error. Coding also allows bits to be reliably received depending on the presence/absence of other users. Hence, collisions are not catastrophic and stability analysis is very different from that of traditional ALOHA.

2.1 Channel Model

The channel model we propose is a multiple access system where M users transmit to one receiver in the presence of additive white Gaussian noise (AWGN). The transmitters all share bandwidth of size W. The signal X_i of user i, along with the output, are bandlimited to W as well. User i is constrained to use an average of up to P_i units of power per transmission. Signals are sampled and synchronized. After sampling, the output and input are related at sample time t as

$$y[t] = \sum_{k=1}^{M} x[t] + N[t]$$

where N[t] is a sequence of $\mathcal{N}(0, \sigma_N^2)$ i.i.d. random variables. Fading is not present in this model.

A user's queue contains all of its traffic which has not yet been successfully transmitted. Immediately before time slot n, $\vec{\Delta}(n) = (\Delta_1(n), \Delta_2(n), ..., \Delta_M(n))$ bits enter the users' queues. During time slot n, $\vec{u}(n) = (u_1(n), u_2(n), ..., u_M(n))$ bits are reliably transmitted and removed from the users' queues. If we denote the number of bits in the users' queues at the beginning of time slot n (after the new bits enter and before transmission begins) as $\vec{Q}(n) = (Q_1(n), Q_2(n), ..., Q_M(n))$, then user *i*'s buffer state evolves according to

$$Q_i(n+1) = \Delta_i(n+1) + Q_i(n) - u_i(n).$$

Once a user receives a packet for transmission, the data in that packet enters the transmission queue and a portion of the data from the queue is transmitted according to a certain policy. Note that some bits may undergo unreliable transmission and need to be retransmitted. Thus the queue holds all bits that need reliable transmission.

Packets arrive to each user's transmission buffers according to independent Bernoulli processes. At each time slot, user *i* has a probability p_i of new packet arrival. Packets are fixed to be L_i . Hence, for each user, the pair (p_i, L_i) characterizes the burstiness of the packet arrivals for user *i*. In the *n*th time slot:

$$\Delta_i(n) = \begin{cases} L_i & \text{with probability } p_i \\ 0 & \text{with probability } 1 - p_i. \end{cases}$$

We assume each user has a buffer of infinite size. Time slots are of length T transmissions. The vector of arrival rates is $\vec{\lambda} = (\lambda_1, \lambda_2, ..., \lambda_M)$ where $\lambda_i = \frac{p_i L_i}{T}$ is the arrival rate (in bits per transmission) to the buffer of user i.

2.2 Data Transmission Policies and State Information

At each time slot, user *i* must either transmit or not transmit over the time slot using codes of length *T*. A collision occurs if two or more receivers transmit during the same time slot. We assume that the receiver and transmitter have perfect synchronization. The receiver knows for each user, at each time slot, whether or not that user is transmitting. This may done by using coded tags on transmissions to identify users. The code on the tags is sufficient to withstand multiple access interference from all users at once. We assume that by the end of each time slot, each user knows which portion of its transmission data has been reliably received, and which portion needs to be retransmitted. We constrain the set of transmission policies of interest at time *n* to be a function of $\vec{Q}(n)$. We may view this as a discrete time controlled stochastic system.

2.3 Markovity and Average Bit Delay

We note that since the buffer arrival process is Bernoulli and our transmission policy at time n is a function of $\vec{Q}(n)$, the system of buffer state evolution satisfies the Markov condition:

$$P[\vec{Q}(n+1) = \vec{q} \mid \vec{Q}(n), \vec{Q}(n-1), ..., \vec{Q}(0)] = P[\vec{Q}(n+1) = \vec{q} \mid \vec{Q}(n)].$$

 \vec{Q} is the state variable of the homogeneous Markov chain. The transition probabilities of the state evolution are governed by the arrival processes burstiness pairs and the transmission policy.

2.4 Coding Time Slots and Collisions

As in traditional time-slotted ALOHA systems, if two or more users attempt to transmit during the same time slot, a *collision* occurs. However, because of our use of coding, a collision is not necessarily catastrophic: data may still be reliably received in the event of other users transmitting. Time slots are of length T transmissions, where T is long enough in terms of bits so that data may be transmitted with acceptable probability of error even in the event of a collision. For a large but finite T, error exponents [9] for multiple access channels [7] quantify at what rate the probability of error decays exponentially with T. Let n be the number of time slots during which we transmit. Each user has K sets of codewords $\mathcal{M}_{j}^{i,k}$, k = 1, ..., K. During time slot j user i has a codebook C_{j}^{i} of length T to encode at most one codeword from each $\mathcal{M}_{j}^{i,k}$. We denote the codebook C_{j}^{i} as $(T, \xi, \lambda_{j}^{i})$ time-slot capacity-achieving if there exists a decoding policy such that $\forall k \in \mathcal{K}_{j}^{\prime i} \subset \{1, ..., K_{j}^{i}\}$: a codeword from $\mathcal{M}_{j}^{i,k}$ was encoded, that codeword was decoded with probability of error ξ or less, and $\sum_{k \in \mathcal{K}_{j}^{\prime i}} \frac{log(|\mathcal{M}_{j}^{i,k}|)}{T} \ge \lambda_{j}^{i}$. Over the long term, for a given T we say that a coding and decoding policy is $(T, \xi, \vec{\lambda})$ *capacity-achieving* if

$$\lim_{n \to \infty} \frac{1}{Tn} \sum_{j=1}^{n} E\Big[\sum_{k \in \mathcal{K}_{j}^{i}} log(|\mathcal{M}_{j}^{i,k}|)\Big] = \lambda_{i}, \qquad 1 \le i \le M,$$

and the probability of error for the transmission of each user on each time-slot is upper bounded by ξ . The notion of capacity described above is related to other delayconstrained and probability-of-failure notions of capacity, such as delay-limited capcity [10], ϵ -capacity [11], capacity versus outage [12, 13], and expected capacity [14]. We consider delay constraints because of finite time slot length. We use expected rates because of uncertainty regarding collisions. The capacity region of our above definition for users with power constraints and certain transmission policy constraints has been found in [5] to be the same as the Cover-Wyner region for multiple access channels.

3 Minimizing Delay and Minizing Power Consumption

Before we proceed to introduce an optimal coding strategy, we attempt to understand how delay and power minimization are affected by knowledge of queue information. We restrict our attention to a two-user scenario, but the results may easily be extended for many users. At the beginning of time slot n,

$$\vec{\bigtriangleup}(n) = \begin{cases} (L_1, L_2) & \text{with probability } p_1 p_2 \\ (L_1, 0) & \text{with probability } p_1 (1 - p_2) \\ (0, L_2) & \text{with probability } (1 - p_1) p_2 \\ (0, 0) & \text{with probability } (1 - p_1) (1 - p_2) \end{cases}$$

3.1 Delay Minimization

We now would like to understand what the minimum amount of power is required to minize delay. Consider the situation where the set of (p_i, L_i) , burstiness pairs, for i = 1, 2, is known to both users. We consider the case where users have full knowledge of each other's queues and when they only have local queue information.

Let us denote $C_{\sigma_N^2}(x) = \frac{1}{2}\log(1 + \frac{x}{\sigma_N^2})$ as the capacity (in bits per transmission) of a discrete-time memoryless Gaussian channel with noise variance σ_N^2 and average power per transmission constraint x. It is the maximum rate at which information may be transmitted with arbitrarily vanishing error probability. Similarly, $C_{\sigma_N^2}^{-1}(x) = \sigma_N^2(2^{2x} - 1)$ is the minimum amount of average power required to transmit rate x and noise variance σ_N^2 with arbitrarily vanishing error probability.

We now comment on the capacity region for the discrete-time two-user AWGN multiple access channel with power constraints $\vec{P} = (P_1, P_2)$ and noise variance σ_N^2 . It is defined to be the subset of \mathbb{R}^2_+ with rate pairs (R_1, R_2) satisfying

$$\sum_{i \in \mathcal{S}} R_i \le C_{\sigma_N^2} \big(\sum_{i \in \mathcal{S}} P_i \big), \qquad \mathcal{S} \subseteq \{1, 2\}.$$

Let us denote the *dominant face* of the capacity region as the subset of all rate pairs in the capacity region that satisfy $R_1 + R_2 = C_{\sigma_N^2}(P_1 + P_2)$. For all rate pairs (R_1, R_2) that are not dominant, there exists a (R_{dom1}, R_{dom2}) that satisfies $R_{dom1} \ge R_1$ and $R_{dom2} \ge R_2$. We note that any dominant rate pair delivers the maximum aggregate rate of reliable transmission for a given power constraint \vec{P} . Or alternatively, the power vector \vec{P} delivers the minimum aggregate power for two users to transmit reliably at rates on the dominant face of the capacity region.

3.1.1 Full Knowledge of Other Users' Queues

Immediately before time slot n, users have full knowledge of the number of bits that have just entered everyone's queue: $\vec{\Delta}(n)$. To minimize delay to 0, each user must empty the total contents of everyone's queues each time slot of length T transmissions. Every user has access to two codebooks: a multiple-access codebook that may achieve the rate pair $(\frac{L_1}{T}, \frac{L_2}{T})$ reliably, and a single-user codebook that may achieve the rate $\frac{L_i}{T}$ reliably for user i when no multiple access interference is present. We note that the minimum amount of aggregate power per transmission required is i.i.d. over each time slot n, and is a function of $\vec{\Delta}(n)$. The minimum amount of aggregate power per transmission required to empty the buffers in one time slot is given by:

$$P_{min}^{(1)}(n) = \begin{cases} C_{\sigma_N^2}^{-1} \left(\frac{L_1 + L_2}{T}\right) & \text{with probability } p_1 p_2 \\ C_{\sigma_N^2}^{-1} \left(\frac{L_1}{T}\right) & \text{with probability } p_1 (1 - p_2) \\ C_{\sigma_N^2}^{-1} \left(\frac{L_2}{T}\right) & \text{with probability } (1 - p_1) p_2 \\ 0 & \text{with probability } (1 - p_1)(1 - p_2) \end{cases}$$

Since the arrival processes are independent and Bernoulli, ergodicity holds and we have:

$$\lim_{m \to \infty} \frac{1}{m} \sum_{n=1}^{m} P_{\min}^{(1)}(n) \xrightarrow[a.s.]{} p_1 p_2 C_{\sigma_N^2}^{-1} \left(\frac{L_1 + L_2}{T}\right) + p_1 (1 - p_2) C_{\sigma_N^2}^{-1} \left(\frac{L_1}{T}\right) + (1 - p_1) p_2 C_{\sigma_N^2}^{-1} \left(\frac{L_2}{T}\right).$$

We note that similar results hold for any set of ergodic arrival processes \triangle .

3.1.2 No Knowledge of Other Users' Queues

If we now assume that each user still has access to the (p_i, L_i) burstiness pairs of everyone, but does not have access to the amount of data entering the other's queue, then to minimize delay, all users must coordinate to transmit at the worst case scenario: when (L_1, L_2) enters each others' queues. So user each only has access to a multiple-access codebook. Each user always anticipates the other user's presence and uses the amount of power required to empty both queues. Note that many power constraints for users may result in the rate pair lying on the dominant face of the multiple access region. Depending on the burstiness probabilities, however, some of these power constraints may provide smaller average power consumption than others. So these power constraints may be chosen as a function of the burstiness pairs, (p_i, L_i) for i = 1, 2, so long as they may reliably achieve the rate pair $(\frac{L_1}{T}, \frac{L_2}{T})$. So if the power constraints lie in the region \mathcal{P} denoted as:

$$P_{1} \geq C_{\sigma_{N}^{2}}^{-1} \left(\frac{L_{1}}{T}\right) P_{2} \geq C_{\sigma_{N}^{2}}^{-1} \left(\frac{L_{2}}{T}\right) P_{1} + P_{2} = C_{\sigma_{N}^{2}}^{-1} \left(\frac{L_{1}+L_{2}}{T}\right)$$

then for block coding multiple access schemes used over a slot, there exists a time-sharing ratio γ such that

$$\gamma C_{\sigma_N^2}(P_1) + (1 - \gamma) C_{\sigma_N^2 + P_2}(P_1) = \frac{L_1}{T} \gamma C_{\sigma_N^2}(P_2) + (1 - \gamma) C_{\sigma_N^2 + P_1}(P_2) = \frac{L_2}{T}$$

The choice of power constraints is made to minimize the long term average aggregate power consumption

$$J(P_1, P_2) = p_1(1-p_2)P_1 + p_2(1-p_1)P_2 + p_1p_2(P_1+P_2) = p_1(1-p_2)P_1 + p_2(1-p_1)\left(C_{\sigma_N^2}^{-1}\left(\frac{L_1+L_2}{T}\right) - P_1\right) + p_1p_2C_{\sigma_N^2}^{-1}\left(\frac{L_1+L_2}{T}\right).$$

So this is in fact a linear objective function in one variable, and must be minimized subject to a constraint on the value P_1 :

$$C_{\sigma_N^2}^{-1}\left(\frac{L_1}{T}\right) \le P_1 \le C_{\sigma_N^2}^{-1}\left(\frac{L_1+L_2}{T}\right) - C_{\sigma_N^2}^{-1}\left(\frac{L_2}{T}\right).$$

The minimum is attained at either of the two boundary points, depending on the sign of $p_1 - p_2$:

$$(P_1^*, P_2^*) = \begin{cases} \left(C_{\sigma_N^2}^{-1} \left(\frac{L_1}{T} \right), C_{\sigma_N^2}^{-1} \left(\frac{L_1 + L_2}{T} \right) - C_{\sigma_N^2}^{-1} \left(\frac{L_1}{\sigma_N^2} \right) \right) & \text{if } p_1 > p_2 \\ \text{any } (P_1, P_2) \in \mathcal{P} & \text{if } p_1 = p_2 \\ \left(C_{\sigma_N^2}^{-1} \left(\frac{L_1 + L_2}{T} \right) - C_{\sigma_N^2}^{-1} \left(\frac{L_2}{T} \right), C_{\sigma_N^2}^{-1} \left(\frac{L_2}{T} \right) \right) & \text{if } p_1 < p_2 \end{cases}$$

So to minimize long-term average power consumption, the rate pair to be achieved lies on either of the two boundary points of the dominant face of the multiple access region for unequal burstiness probabilities. The long term average minimum amount of power per transmission required for this scheme is given by

$$\lim_{m \to \infty} \frac{1}{m} \sum_{n=1}^{m} P_{\min}^{(2)}(n) \xrightarrow[a.s.]{} p_1(1-p_2)P_1^* + p_2(1-p_1)P_2^* + p_1p_2(P_1^*+P_2^*),$$

where P_1^* and P_2^* have been given above.

3.2 Minimizing Power Consumption

We now address the minimum amount of average power consumption needed to stabilize the bursty system. Concavity of the function $\log(1 + x)$ provides the inequality $\frac{1}{2}\log(1+\frac{P}{\sigma_N^2}) \ge 2 * \frac{1}{2}\log\left(1+\frac{P}{\sigma_N^2}\right)$. So it is more favorable in terms of aggregate power consumption to spread the same amount of power into multiple time slot uses rather than in just one use. Note that we may generalize this to more than two slots, so long as the system is stable. For multiple users to reliably transmit at a prescribed rate-tuple, using a coding scheme where that rate-tuple lies on the dominant face of the multiple access capacity region minimizes the amount of aggregate power required. Note that in terms of long-term power consumption, for certain types of ergodic arrival processes, user queue information is not necessary to perform this strategy. If each user artificially backs up its queue by not transmitting, then after a while each user will have very large queue lengths. At that point, each user will have data to transmit. Afterwards, users transmit achieving the rate pair $(\frac{p_1L_1}{T}, \frac{p_2L_2}{T})$ lying on the dominant face of the Cover-Wyner region. As the vector of output rates tends toward the vector of input rates from above, the amount of power consumption required decreases, but average delay increases. Since the system must provide stability, the minimum amount of power required will correspond to when the vector of output rates matches the vector of input rates with equality. For ergodic processes, the proportion of time users spent artificially backing up queues tends to 0. We note, that since each of the server utilizations is exactly 1, the average delay is infinite. The average aggregate amount of power required is given by

$$\lim_{m \to \infty} \frac{1}{m} \sum_{n=1}^{m} P_{min}^{(3)}(n) \xrightarrow[a.s.]{} C_{\sigma_N^2}^{-1} \left(\frac{p_1 L_1}{T} + \frac{p_2 L_2}{T}\right) = C_{\sigma_N^2}^{-1} (\lambda_1 + \lambda_2).$$

4 Power and Delay Tradeoffs for a System with Limited Queue Information

We now present a coding scheme that addresses the burstiness of packet arrivals and analyze how it performs. The system utilizes a very small amount of queue information among users to operate in two modes: a multiple-access mode when all queue lengths are large, and a broadcast mode otherwise. This scheme tries to address both delay and aggregate power consumption parameters by affording a compromise between the schemes mentioned in the previous section.

4.1 System Design

We assume a limited information sharing scheme where, at time slot n, each user does not know the contents of the newly arrived packets in other users queues. By the end of the time slot, perhaps through feedback from the receiver, each user knows which portion of the data it attempted to transmit was received reliably, and which portion needs to be retransmitted.

4.1.1 Large Queue Lengths: Multiple Access Mode

The results in [5] show that the capacity region of the time-slotted ALOHA system with power-constrained users is the same as the capacity region of the multiple-access channel. For any vector of arrival rates lying inside the Cover-Wyner region, there exist coding schemes that will provide system stability. As described in [5], as all users' queues become very backed up, they are able to transmit simultaneously at rate-tuples lying on the dominant face of the multiple access region while sustaining a small upper bound on probability of error. Error exponents we discussed previously provide bounds to the probability of error for a given slot length, T.

For our system Markov chain with state variable \vec{Q} , we denote the vector-valued drift of the state to be

$$\vec{D}(\vec{q}) = E[\vec{Q}(n+1) - \vec{Q}(n) \mid \vec{Q}(n) = \vec{q}].$$

In our case, based on our model of the data transmission policies state information in Section 2.2, $D_i(\vec{q}) = \lambda_i - \mu_i(\vec{q})$, where $\mu_i(\vec{q})$ is a function of \vec{q} . If $D_i(\vec{q}) < 0 \ \forall i$, then via Pak's Lemma [15], the chain is ergodic, and steady-state probabilities exist. Hence, a sufficient condition for stability of our system model is to is to eventually transmit at multiple access after the queue states cross a finite threshold $\vec{\eta} = (\eta_1, \eta_2, ..., \eta_M)$. Thus, given a set of burstiness pairs and per transmission power constraints, as all users' queues states become backed up (cross this threshold $\vec{\eta}$), they transmit data out of the queues at rates $\vec{\mu}_{ma} = (\mu_{ma1}, \mu_{ma2}, ..., \mu_{maM}) = E[\vec{\mu} \mid \vec{Q} \geq \vec{\eta}]$. To transmit optimally aggregate data subject to the power constraints, $\vec{\mu}_{ma}$ must lay on the dominant face of the multiple access region: $(\sum_j \mu_{maj} = C_{\sigma_N^2}(\sum_j P_j))$. To ensure stability that operating point should provides negative drift for all queues $(\lambda_j - \mu_{maj} < 0 \ \forall j)$. We note that the particular threshold is given by $\eta_j = T\mu_{maj}$, where T is the time slot length in transmissions.

4.1.2 Small Queue Lengths: Broadcast Mode

When all users' queues lengths are not above the threshold $\vec{\eta}$, they switch to a broadcast mode where they may combat burstiness by achieving variable reliably received rates. Capacity on the degraded AWGN broadcast channel is achieved by rate-splitting (where a user superimposes two independent virtual-user codes) at the transmitter and successive decoding (signals are iteratively decoded and subtracted out for future decoding). One (low-resolution) signal is coded to be received reliably by both users. The other (highresolution) signal is coded to be received reliably by the receiver with the stronger SNR. The receiver decodes a virtual user, eliminates its contribution, and then decodes the next.

In the small queue length regime of our system, we use this broadcast idea to combat burstiness: each user splits into virtual users which code anticipating different users' presence that time slot. Even in the event of a collision, data is reliably received from all users. Each user has a simple *deterministic* transmission policy: if a user has data in its queue to transmit, it attempts to do so. The only time a collision does not take place is when one of the two users' queues is empty. Each user codes to transmit over two possible channels: a channel with the other user present, and a channel without the other user. A fraction α of user *i*'s power P_i , is allocated to a virtual low-resolution user that codes anticipating not only the presence of the virtual user counterpart for user *i*, but also the other physical user's presence. The high-resolution virtual user for user *i* does not anticipate the other physical users' presence: it is only received reliably when user $j \neq i$ does not transmit. The rates at which each virtual user attempts to transmit reliably via successive decoding are as follows:

$$\mu_{LR1} = \frac{1}{2}log(1 + \frac{\alpha_1 P_1}{\sigma_N^2 + P_2 + (1 - \alpha_1)P_1}))$$

$$\mu_{HR1} = \frac{1}{2}log(1 + \frac{(1 - \alpha_1)P_1}{\sigma_N^2}))$$

$$\mu_{LR2} = \frac{1}{2}log(1 + \frac{\alpha_2 P_2}{\sigma_N^2 + P_1 + (1 - \alpha_2)P_2}))$$

$$\mu_{HR2} = \frac{1}{2}log(1 + \frac{(1 - \alpha_2)P_2}{\sigma_N^2}).$$

In each time slot, a collision occurs when more than one user transmits. If a collision occurs, only the low-resolution component of each user is reliably received. Otherwise, the low and high-resolution components of the sole transmitting user are reliably received. Hence, when the system is in this mode, the reliably received rate pair is as follows:

$$\vec{\mu}_{BC} = \begin{cases} (\mu_{LR1}, \mu_{LR2}) & \text{if a collision occurs} \\ (\mu_{LR1} + \mu_{HR1}, 0) & \text{if only user 1 transmits} \\ (0, \mu_{LR2} + \mu_{HR2}) & \text{if only user 2 transmits} \end{cases}$$

4.2 Queue Information Sharing

We note that users have different sets of codebooks for which they transmit information, a set of codebooks for when they transmit in multiple access mode, and a set of codebooks for when they transmit in broadcast mode. Users notify each other when their queue state crosses the threshold η_i . Hence, each user has total knowledge of a synchronized finite-state automaton that denotes whether or not each user's queue length has crossed η_i . When the FSA is in the state where all users thresholds have are greater than η_i , each user switches to multiple access mode. Otherwise, they operate in their broadcast channel mode. We do not model the communication link between users for this communication, but note that it is not a substantial amount of information that is shared.

4.3 Performance

We note that for any set of burstiness pairs $\{(p_i, L_i)\}_{i=1}^2$ with corresponding rate-tuples $\left(\frac{p_1L_1}{T}, \frac{p_ML_2}{T}\right)$ lying inside the Cover-Wyner region, proper coding of our scheme during multiple access mode will result the Markov chain being ergodic. Let us consider the two-user scenario and note how the analysis may easily extended for more users. The steady-state probabilities $\pi_{\vec{q}} = \lim_{n \to \infty} P[\vec{Q}(n) = \vec{q}]$ for the Markov chain are governed by:

- The average per-transmission power constraints P_i for each user,
- The burstiness pairs (p_i, L_i) of each user, and
- The rate-splitting power ratio α_i for each user

We note that the long-term average queue size $N(\vec{P}, \vec{p}, \vec{L}, \vec{\alpha}) = \sum_{\vec{q} \in \mathbb{R}^M_+} \sum q_i \pi(\vec{P}, \vec{p}, \vec{L}, \vec{\alpha})$ may be used to calculate the long-term average bit delay $\overline{T}(\vec{P}, \vec{p}, \vec{L}, \vec{\alpha})$ via Little's Result: $\overline{T} = \frac{N}{\lambda_1 + \lambda_2}$. Since the chain is ergodic for all arrival rates inside the multiple access region, we may truncate the state space and perform an approximation [16] using simulations on a state space of a finite number of states.

Since we use a very limited amount of queue information sharing, and are willing to accept a small but non-zero average bit delay, this proposed scheme affords a compromise between the delay minimizing scheme with no queue information and the power consumption minimizing system of the previous section. We denote the *reasonable* power constraint region as the set of power constraints for our system that satisfy:

$$C_{\sigma_N^2}^{-1} \left(\frac{p_1 L_1}{T} + \frac{p_2 L_2}{T} \right) \le P_1 + P_2 \le C_{\sigma_N^2}^{-1} \left(\frac{L_1 + L_2}{T} \right).$$

If the power constraints were to lie below the lower bound, the system would not be stable $(\vec{\lambda} > \vec{\mu})$. On the other hand, if the power constraints were to lie above the upper



Figure 2: average aggregate power consumption and average bit delay as a function of burstiness for fixed packet length and varying probabilities with $\alpha_1 = \alpha_1 = 0.5$, $\sigma_N^2 = 1$

bound, then the corresponding power constraints of a delay-minimizing scheme would be less, and so would the average delay (which is 0).

Figure 2 shows simulation results for the average power consumption and average bit delay for our proposed scheme with power constraints at the midpoint of the boundaries of the reasonable power constraint region for each value of $p = p_1 = p_2$. The average power consumption of the systems mentioned in the previous section are superimposed in the figure as well. In the regime of small yet nonzero burstiness probabilities (which is where most bursty packetized systems operate), the impact of allowing a small yet nonzero tolerable delay along with a small amount of queue information sharing is illustrated: both average bit delay and average power consumption are near their respective minimal boundaries. Our scheme uses less energy than that of a system with no queue information and 0 delay because that system obtains no large benefit from one of the two users being empty. In our scheme, however, most of the time the system is in broadcast mode and if one of the two users is empty, that user consumes no power to transmit while other user may reliably transmit the low-resolution and high-resolution data during that time slot. Thus, the user's queue length is decreased more rapidly and may tend to 0 faster, where it will not consume power again. As the burstiness probability and henceforth $\dot{\lambda}$ increases, the system spends a majority of time in broadcast mode, and user queue lengths increase along with average delay and power consumption. Eventually, as the arrival rate continues to increase, user queue lengths cross $\vec{\eta}$ more often and the system spends more of its time in multiple access mode, where it transmits at optimal aggregate rates (as given by the dominant face of the Cover-Wyner region). This is illustrated in how the increase in delay slows down. At very high burstiness probabilities, the interval of reasonable power constraints is quite narrow. Hence, the power constraint required to deliver 0 delay is not much larger than the requirement for stably minimizing power consumption, which yields infinite delay. Since the power constraints of our scheme lie within this interval, and the system arrival process is in effect becoming deterministic, the average delay in this high burstiness probability regime becomes arbitrarily large.

References

[1] J. I. Capetanakis, "The multiple access broadcast channel: Protocol and capacity considerations," *MIT Tech. Rep ESL-R-806*, 1978.

- [2] R. Rivest, "Network control by Bayesian broadcast," IEEE Transactions on Information Theory, vol. 33, pp. 323–328, 1987.
- [3] J. Hayes, "An adaptive technique for local distribution," *IEEE Transactions on Communications*, vol. 26, pp. 1178–1186, 1978.
- [4] S. Ghez, S. Verdú, and S. Schwartz, "Stability properties of slotted ALOHA with multipacket reception capability," *IEEE Transactions on Automatic Control*, vol. 33, pp. 640–649, 1988.
- [5] M. Médard, S. P. Meyn, J. Huang, and A. J. Goldsmith, "Capacity of time-slotted ALOHA packetized multiple-access systems," 2001.
- [6] R. A. Berry, Power and Delay Trade-offs in Fading Channels, PhD dissertation, Massachusetts Institute of Technology, Department of Electrical Engineering and Computer Science, June 2000.
- [7] H. Liau, *Multiple Access Channels*, PhD dissertation, University of Hawaii, Department of Electrical Engineering and Computer Science, June 1972.
- [8] T. M. Cover, "Broadcast channels," *IEEE Transactions on Information Theory*, vol. 18, pp. 2–14, 1972.
- [9] R. Gallager, Information Theory and Reliable Communication, John Wiley & Sons, New York, NY, 1968.
- [10] S.V. Hanly and D.N.C. Tse, "Multiaccess fading channels. ii. delay-limited capacities," *IEEE Transactions on Information Theory*, vol. 44, no. 7, pp. 2816 – 2831, 1998.
- [11] S. Verdú and T. S. Han, "A general formula for channel capacity," *IEEE Transac*tions on Information Theory, vol. 40, no. 7, 1994.
- [12] S. Shamai, "A broadcast strategy for the Gaussian slowly fading channel," International Symposium on Information Theory, p. 150, 1997.
- [13] G. Taricco G. Caire and E. Biglieri, "Minimum outage probability for slowly-fading channels," *Proceedings of the International Symposium on Information Theory*, p. 7, 1998.
- [14] M. Effros and A. Goldsmith, "Capacity definitions and coding strategies for general channels with receiver side information," *Proceedings of ISIT*, p. 39, 1998.
- [15] A. G. Pakes, "Some conditions for ergodicity and recurrence of Markov chains," Oper. Res., vol. 17, pp. 1059–1061, 1969.
- [16] D. Freedman, Approximating Countable State Markov Chains, Holden-Day, San Francisco, CA, 1971.